

GPU-accelerated CFD Simulations for Turbomachinery Design Optimization

Mohamed H. Aissa

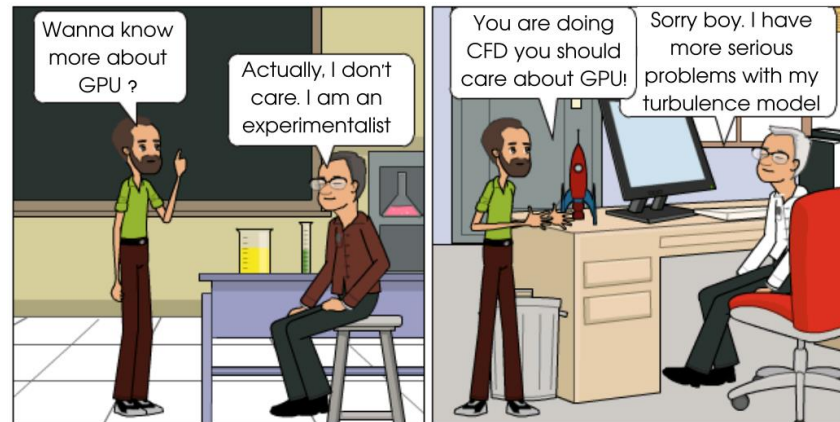
www.researchgate.net/profile/Mohamed_Aissa3

Co-promotor: Dr. Tom Verstraete

Promotor: Prof. C. Vuik



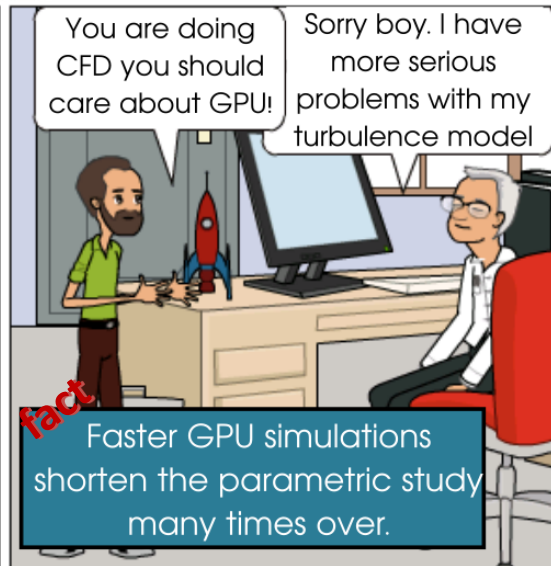
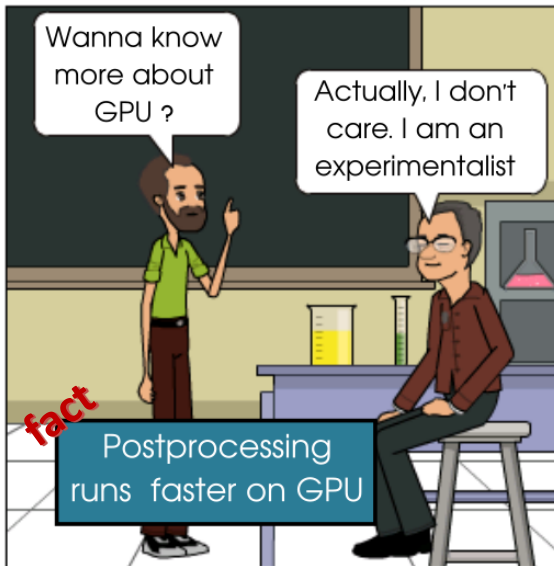
VON KARMAN INSTITUTE
FOR FLUID DYNAMICS



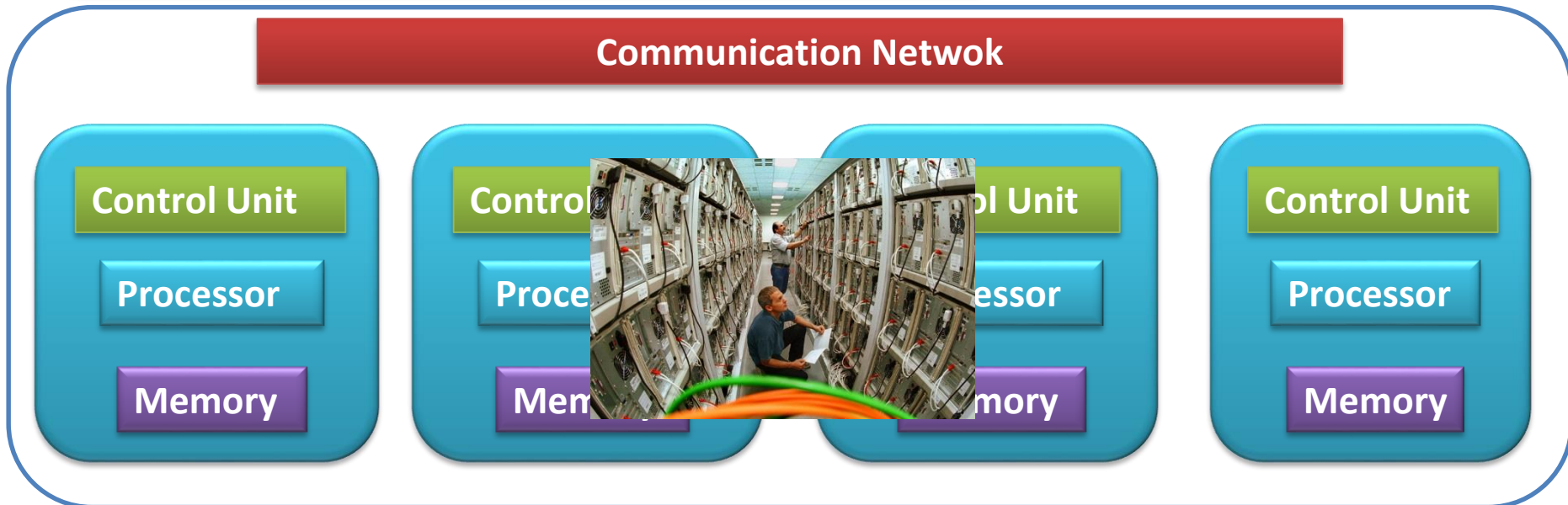
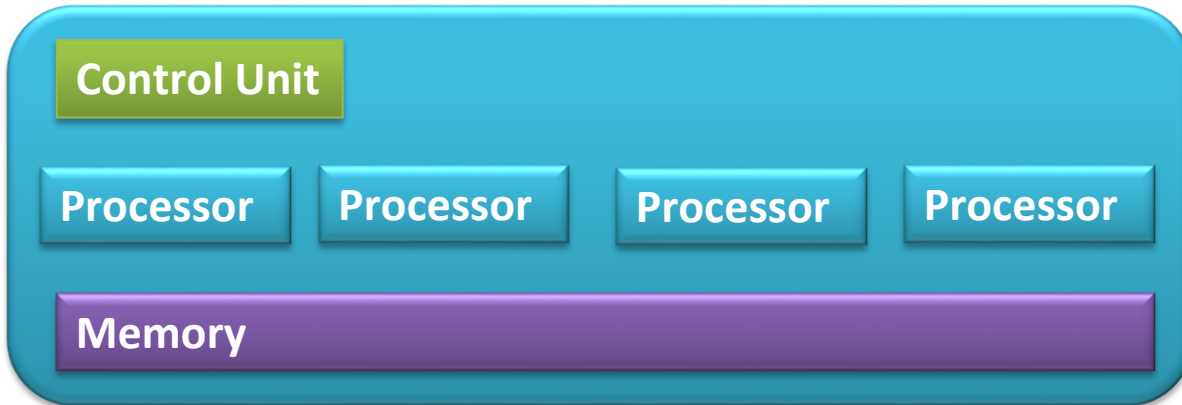
Can your simulation profit from the GPU?



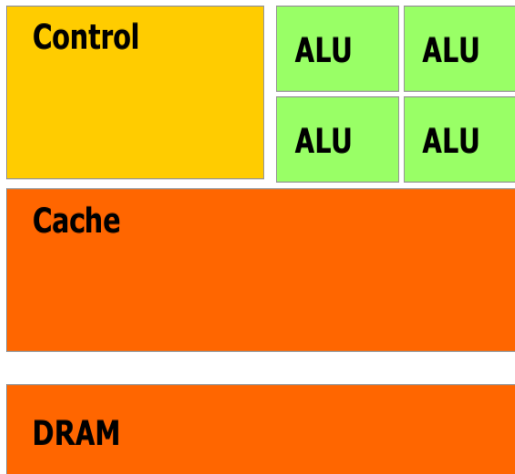
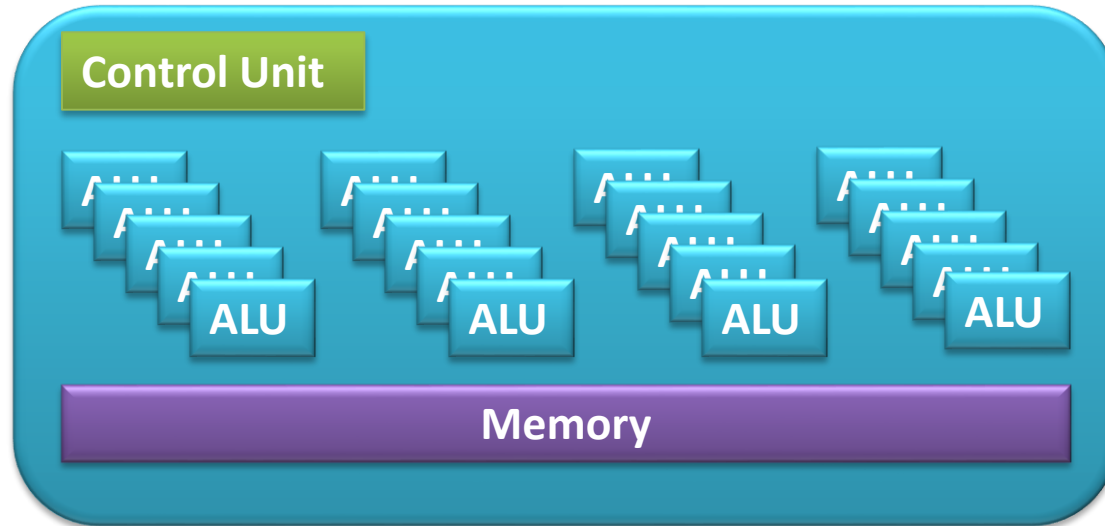
- What is a GPU?
- How fast it is?
- How to use it?



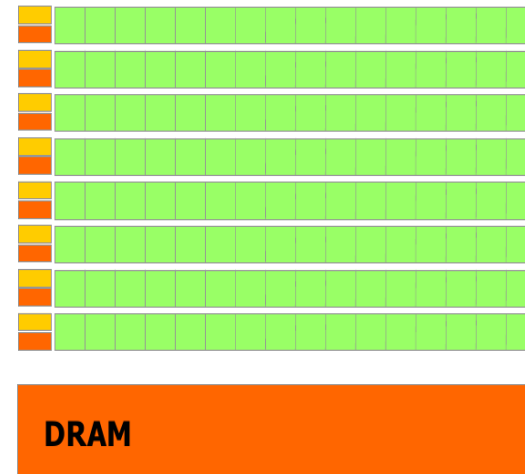
Multi-core vs many-core



Massive Parallel Systems (e.g. GPU) as a trade-off

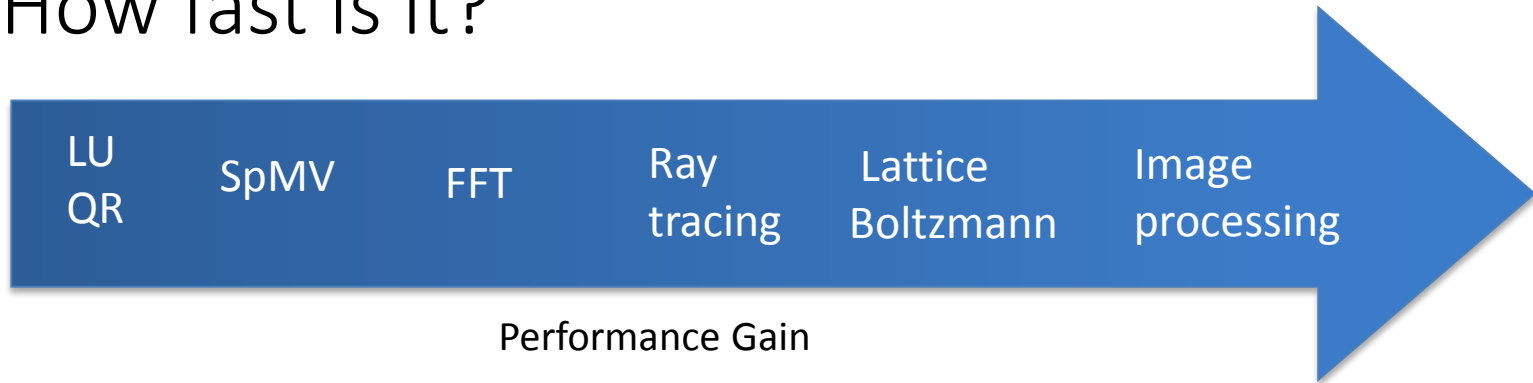


CPU

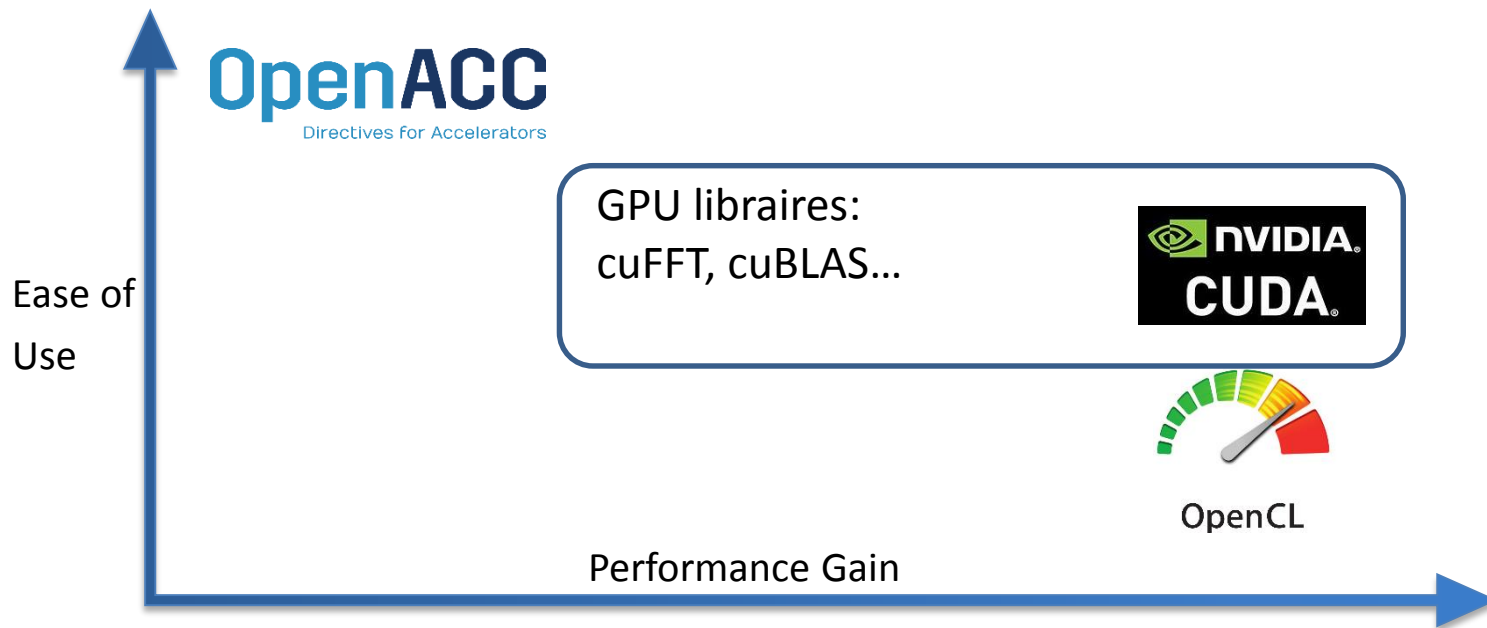


GPU

How fast is it?



How to use a GPU



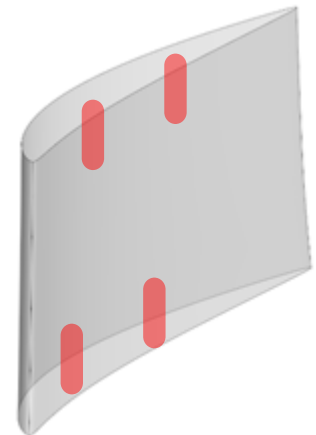
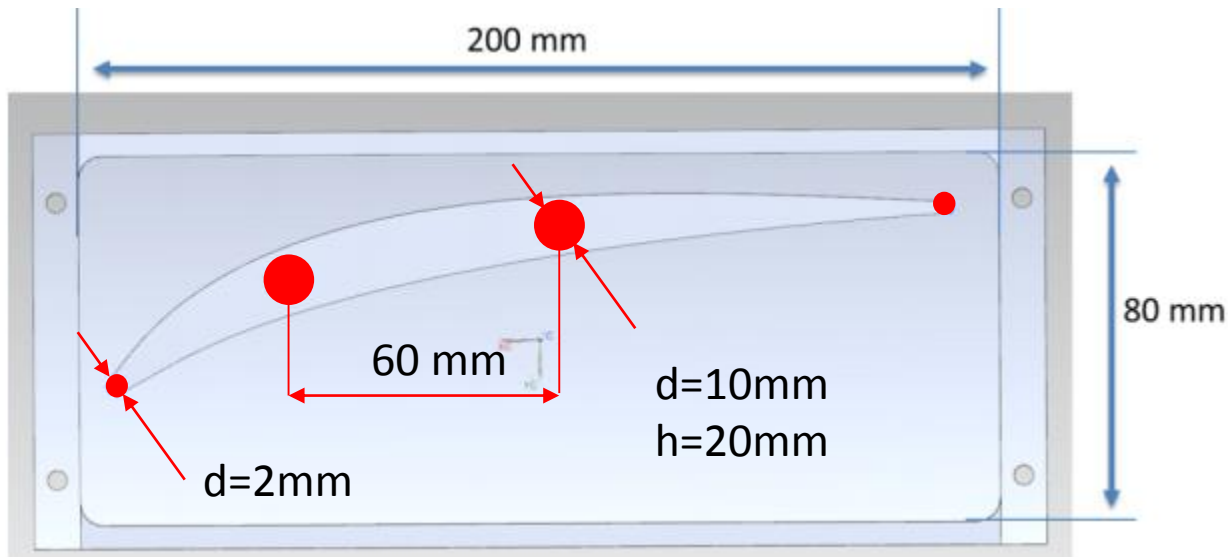
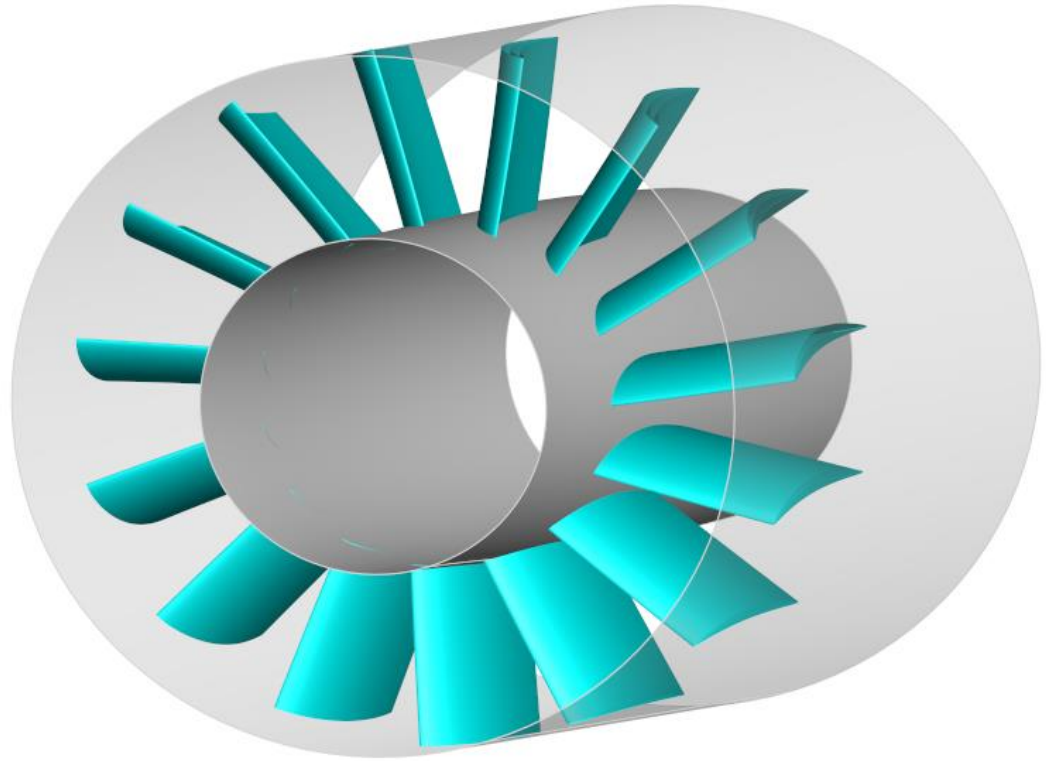
Airplanes are getting
more efficient

Engine optimization
is a main contributor



TurboLab Stator (1/4)

- $N_{\text{blades}} = 15$
- Chord length fixed
- Casing fixture



TurboLab (2/4): Boundary conditions and summary

Inlet P_0 : 102713.0 Pa
Inlet T_0 : 294.314 K

9 kg/s

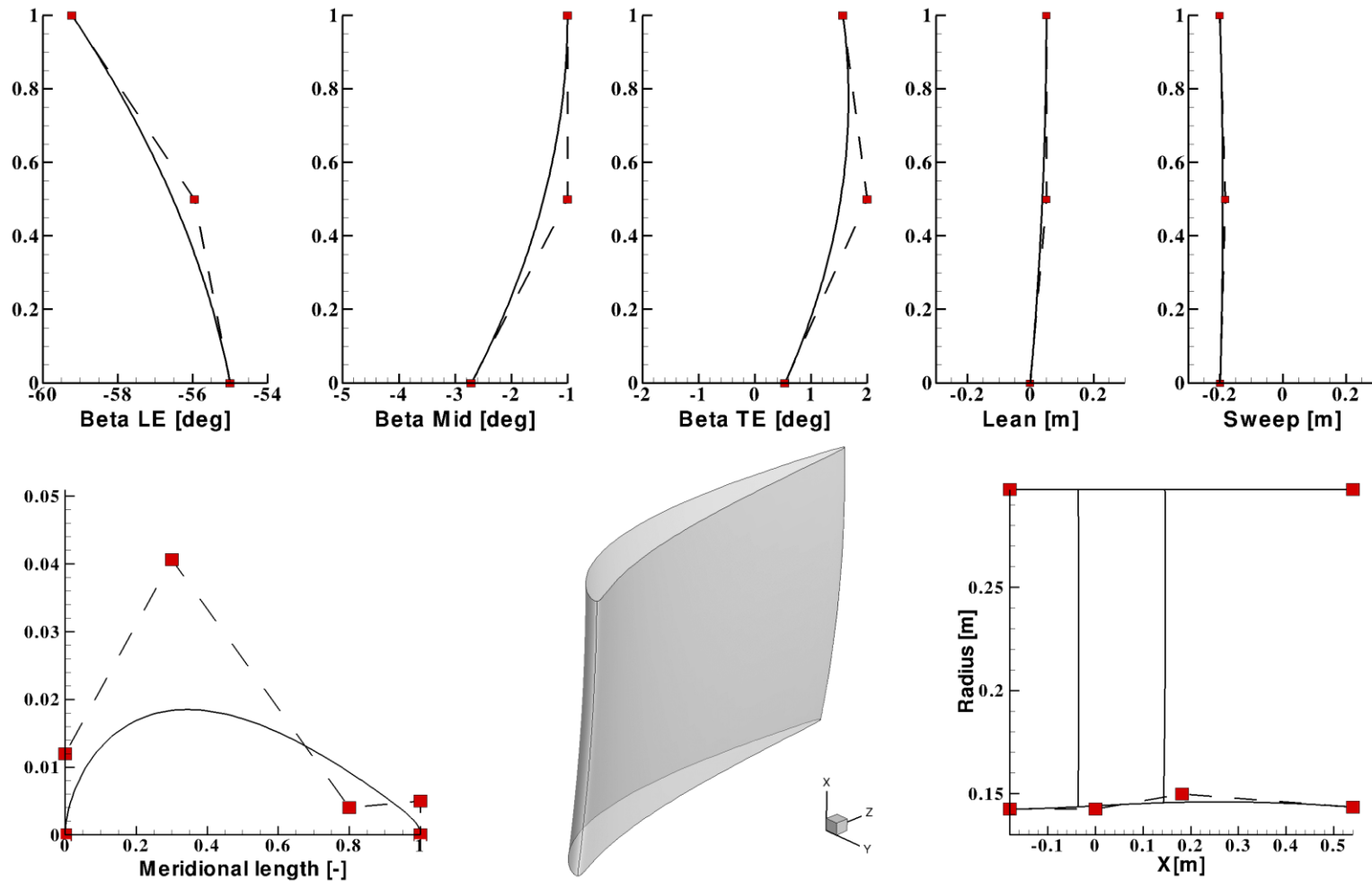
Objectives:

- Lower axial deviation
- Lower total pressure loss

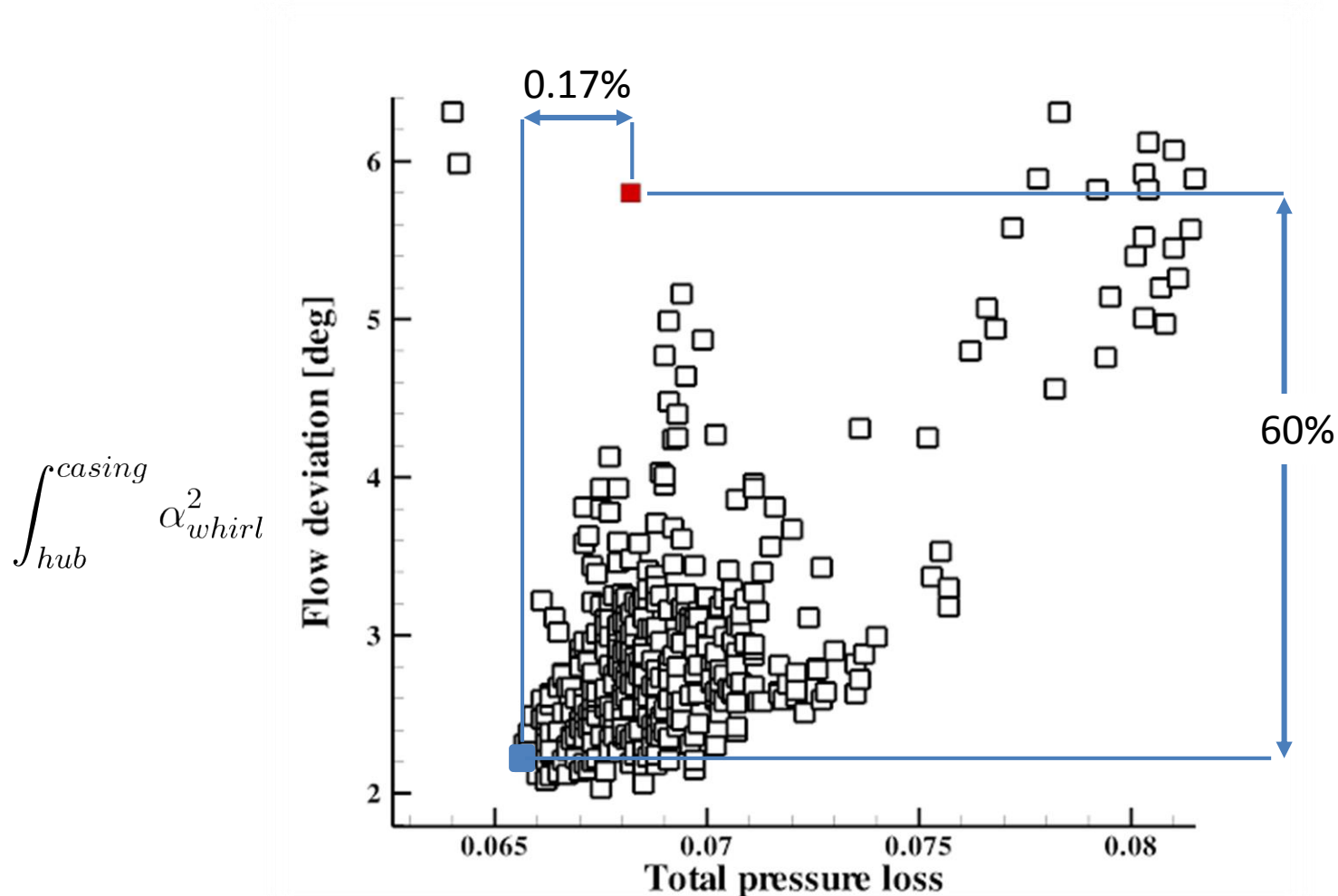
Inlet whirl angle: 42°
Inlet pitch angle: 0°

TurboLab (3/4): Parametrization

21 Design variables



TurboLab (4/4): Optimization Results



$$Loss_{P_0} = \frac{p_{01} - p_{02}}{p_{01} - p_1}$$

Every point is a costly CFD optimization → need for a HPC solution

How beneficial are GPUs, a quick literature check:

- Acceleration is case-dependent (from 1x to 1000x).
- Speedups are sometimes contradicting.
- Some publications are very critical to GPUs for scientific computations:
 - Lee et al “Debunking the 100x GPU vs. CPU myth”
 - Vuduc et al. “On the limits of GPU acceleration”

Main objective:

A more tangible GPU potential



CFD GPU solvers



Classification
of CFD operations



Proof-of-concept:
Optimization cases



Summary
and Conclusions

Main objective:

A more tangible GPU potential



CFD GPU solvers



Classification
of CFD operations



Proof-of-concept:
Optimization cases



Summary
and Conclusions

Numerical Scheme:

$$\frac{\partial}{\partial t} \int_{\Omega} \mathbf{W} d\Omega + \oint_{\partial\Omega} (\mathbf{F}_e - \mathbf{F}_v) dS = \int_{\Omega} \mathbf{Q} d\Omega \quad \mathbf{W} = \{\rho, \rho V_x, \rho V_y, \rho V_z, \rho E\}$$

$$\frac{\Omega}{\Delta t} \Delta \vec{W}^n \left\{ \begin{array}{ll} \text{Explicit Time Stepping} & \Delta \mathbf{W}^n = -\frac{\Delta t}{\Omega} \mathbf{R}^n \\ \text{Implicit Time Stepping} & \left[\frac{(\Omega I)}{\Delta t} + \left(\frac{\delta \mathbf{R}}{\delta \mathbf{W}} \right) \right] \Delta \mathbf{W}^n = -\mathbf{R}^n \end{array} \right.$$

Main objective:

A more tangible GPU potential



CFD GPU solvers

- Explicit time integration



Classification
of CFD operations



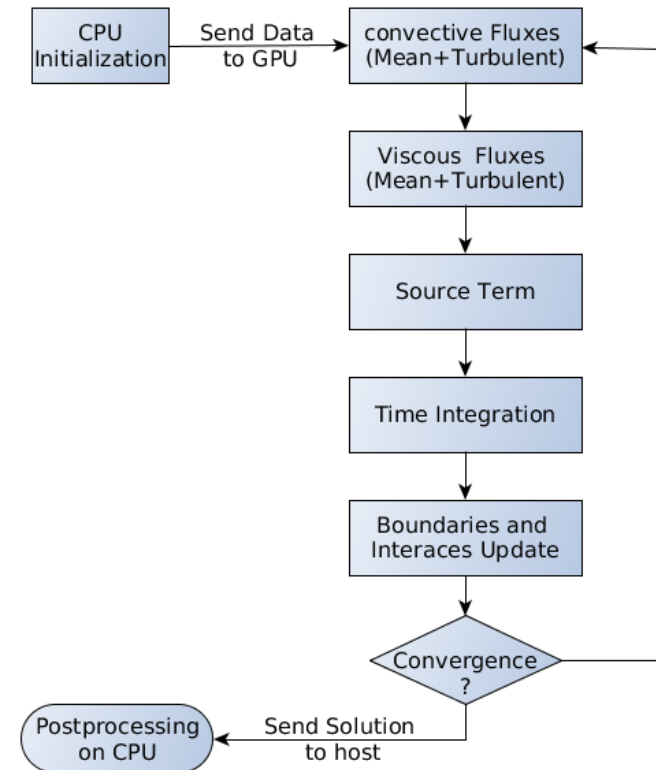
Proof-of-concept:
Optimization cases



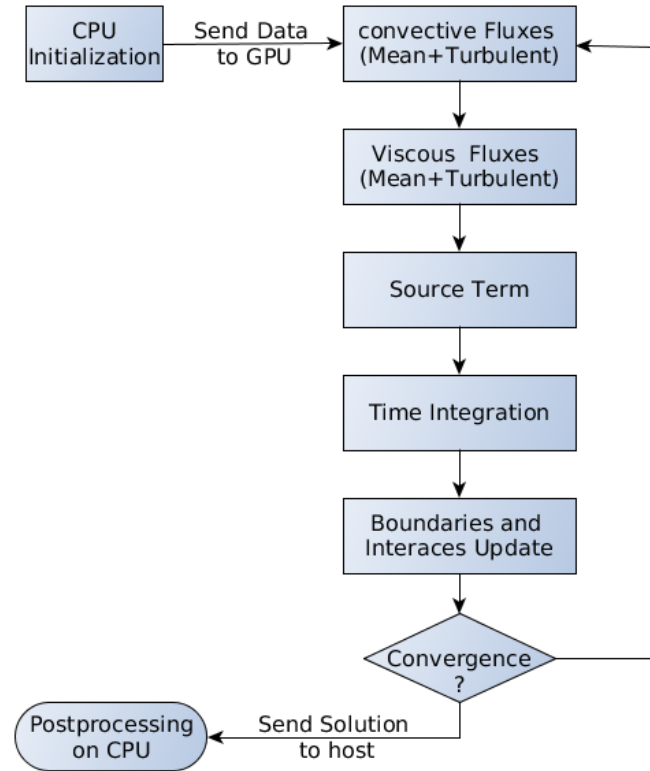
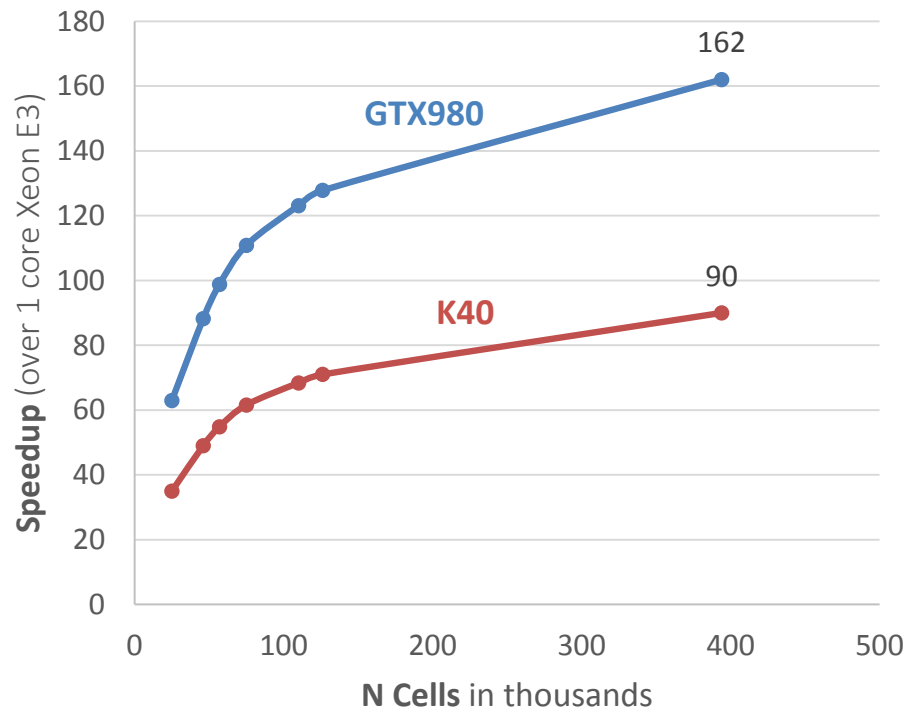
Summary
and Conclusions

Explicit solver

- Application:
Steady RANS simulation
- Solved Equations:
RANS (SA Model)
- Discretization (2nd Order):
Roe Scheme + Flux Limiter
Explicit RK 4 Stage
- Mesh: Multi-Block, Structured
- Acceleration:
 - 2 level Multigrid
 - Implicit Residual Smoothing

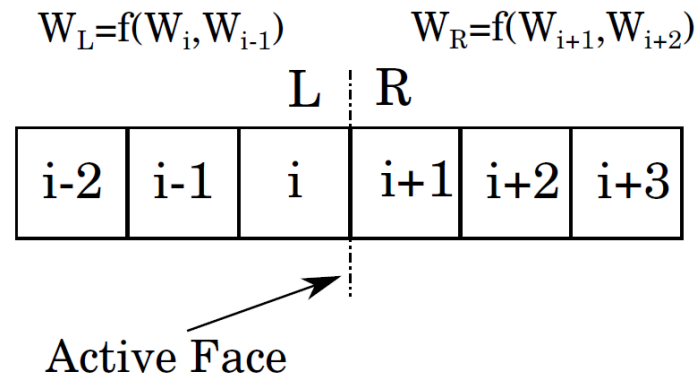


Explicit solver



Convective Flux Evaluation (1/3)

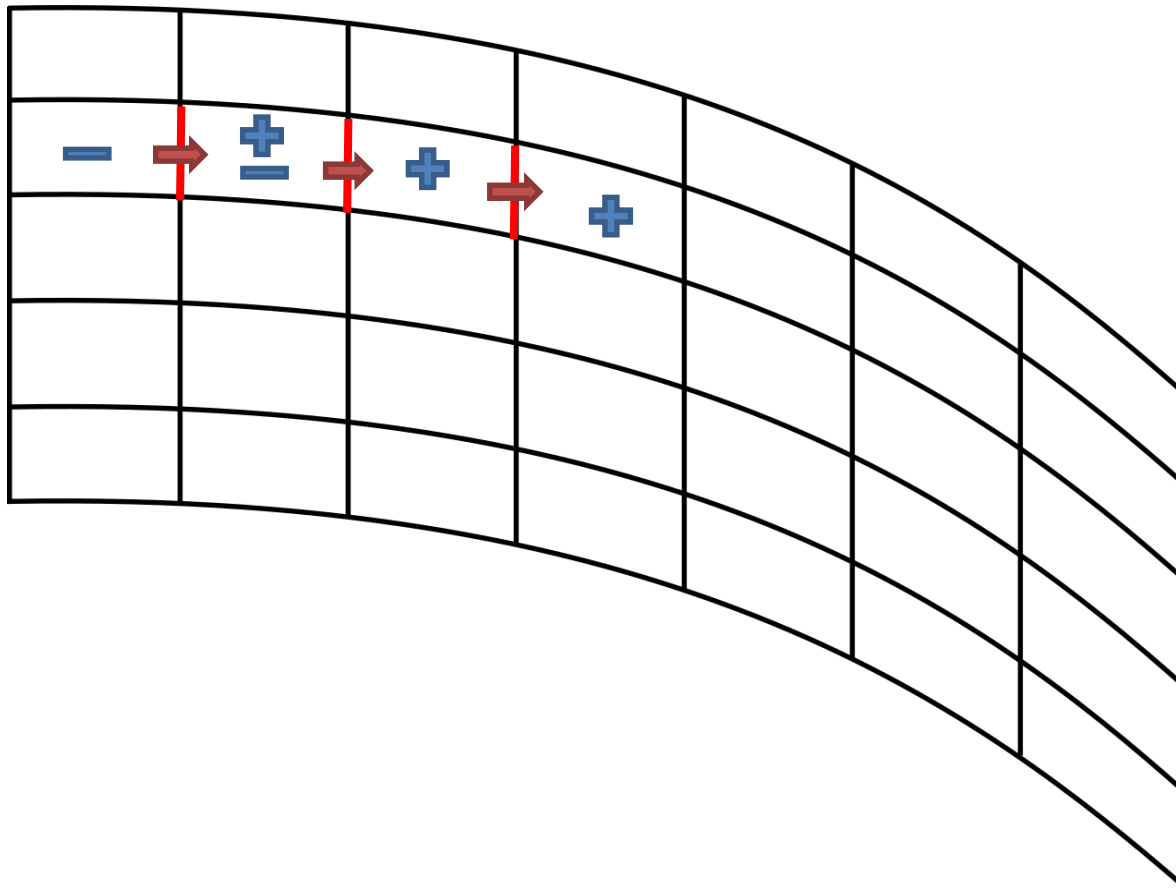
$$(\vec{F}_c)_{I+1/2} = \frac{1}{2}[\vec{F}_c(\vec{W}_R) + \vec{F}_c(\vec{W}_L) - |\bar{A}_{Roe}|_{I+1/2}(\vec{W}_R - \vec{W}_L)]$$



Convective Flux (2/3):

Thread mapping possibilities

- Face-wise is not thread-safe

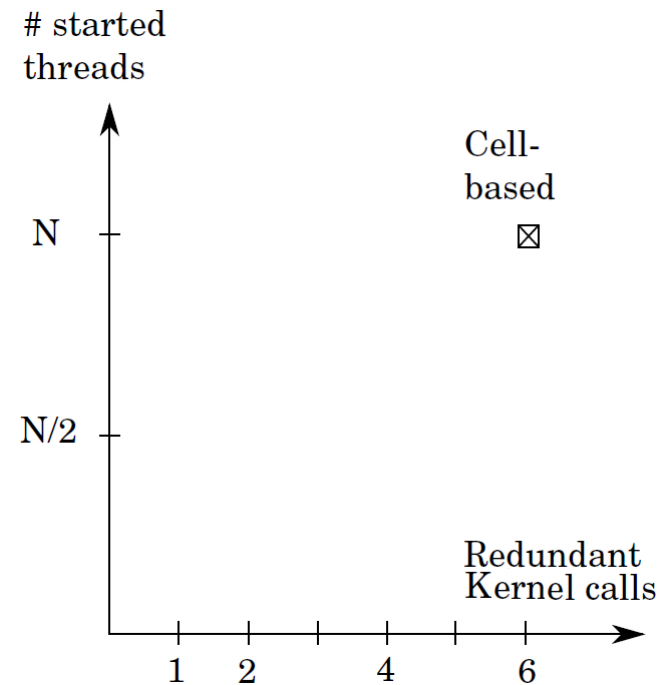
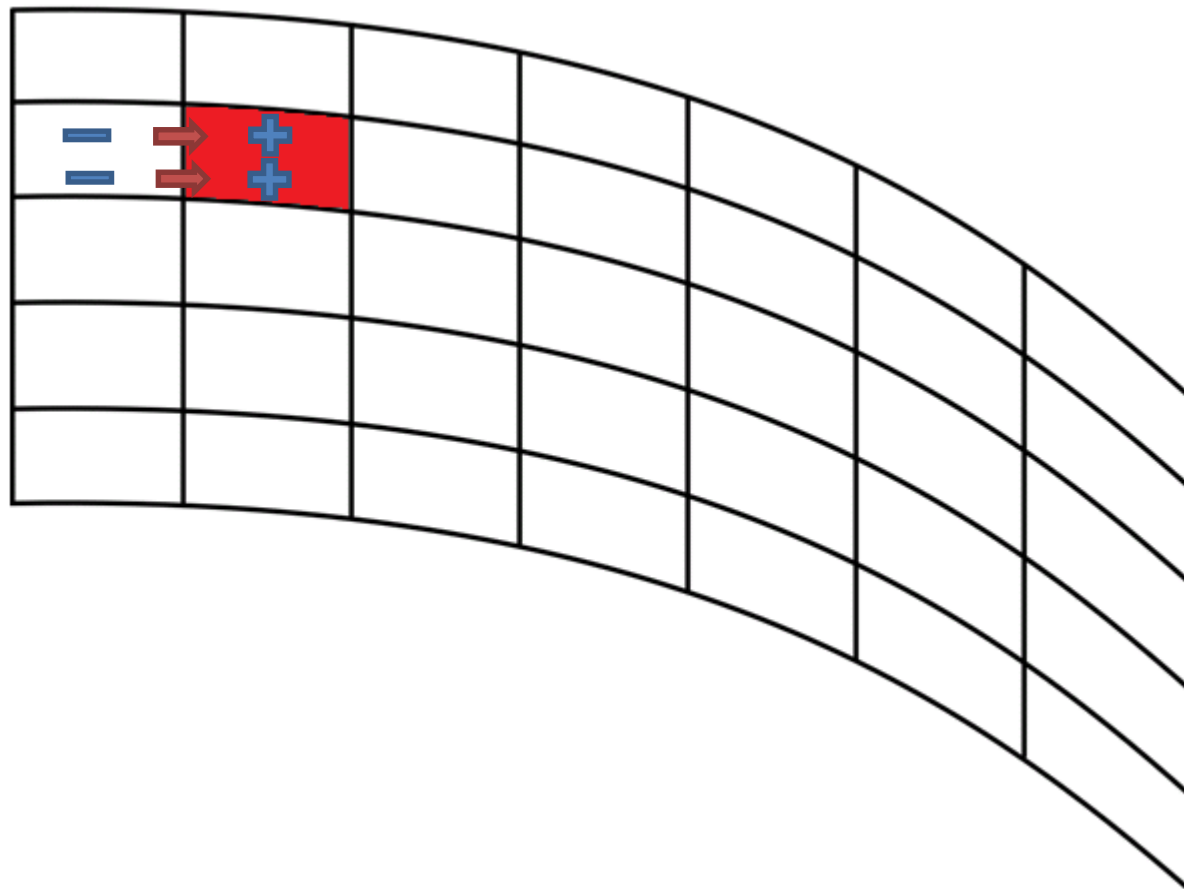


Convective Flux (2/3):

Thread mapping possibilities

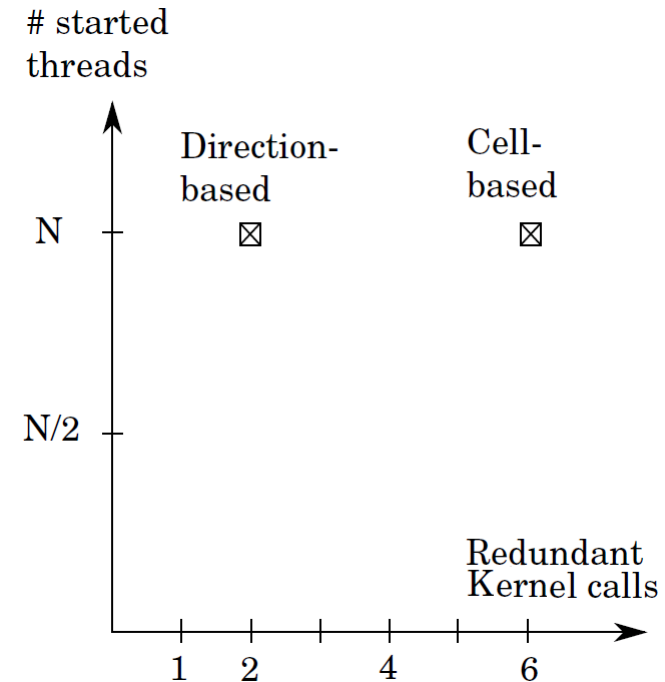
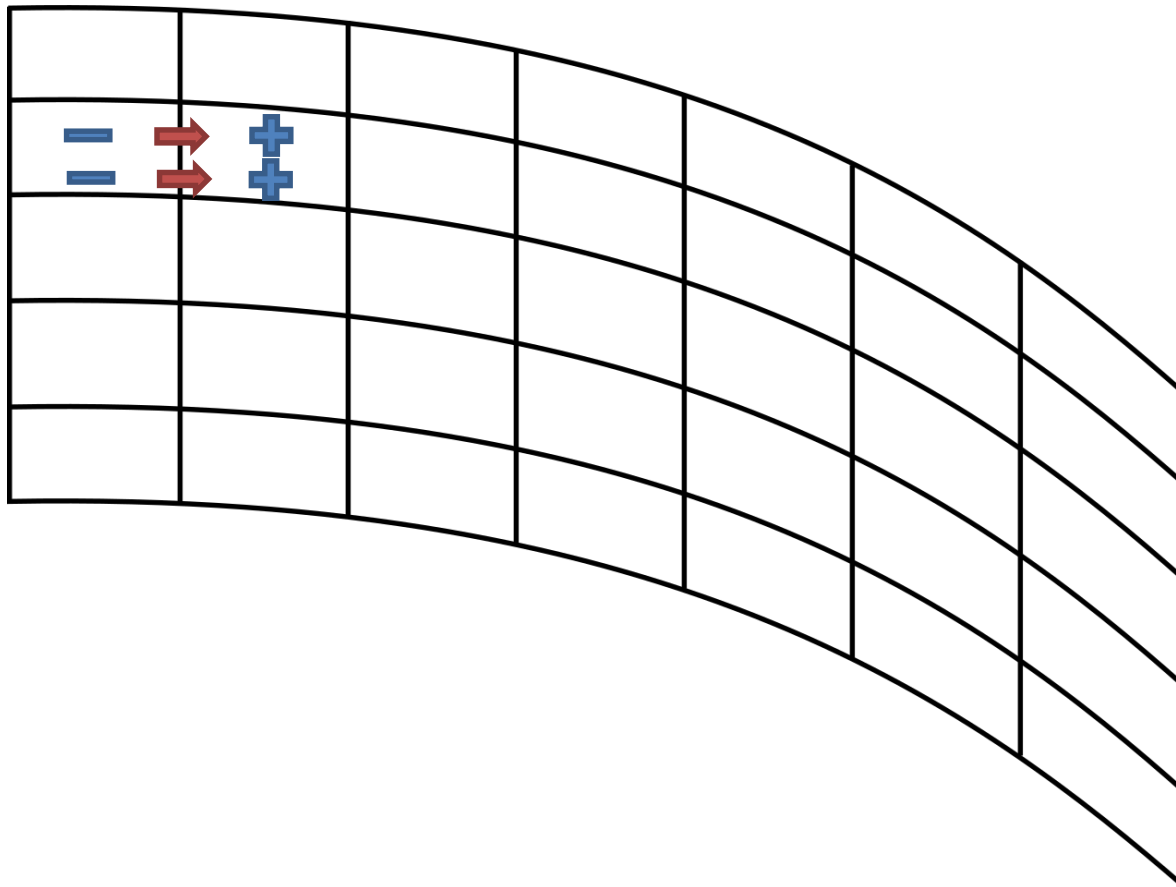
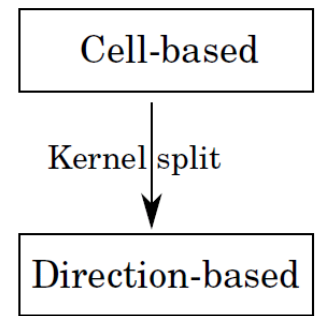
Cell-based

- Cell-based mapping thread safe but with redundancy



Convective Flux (2/3): Thread mapping possibilities

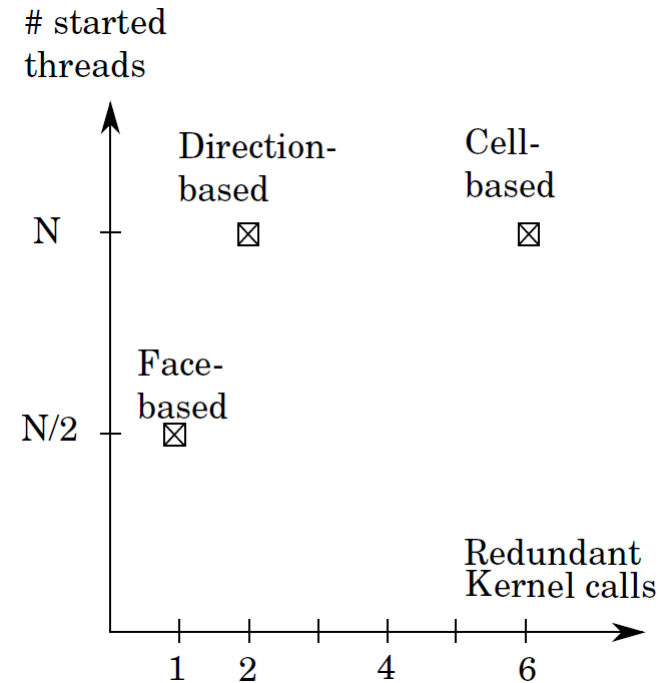
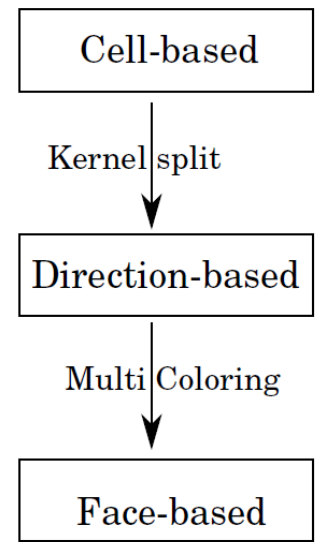
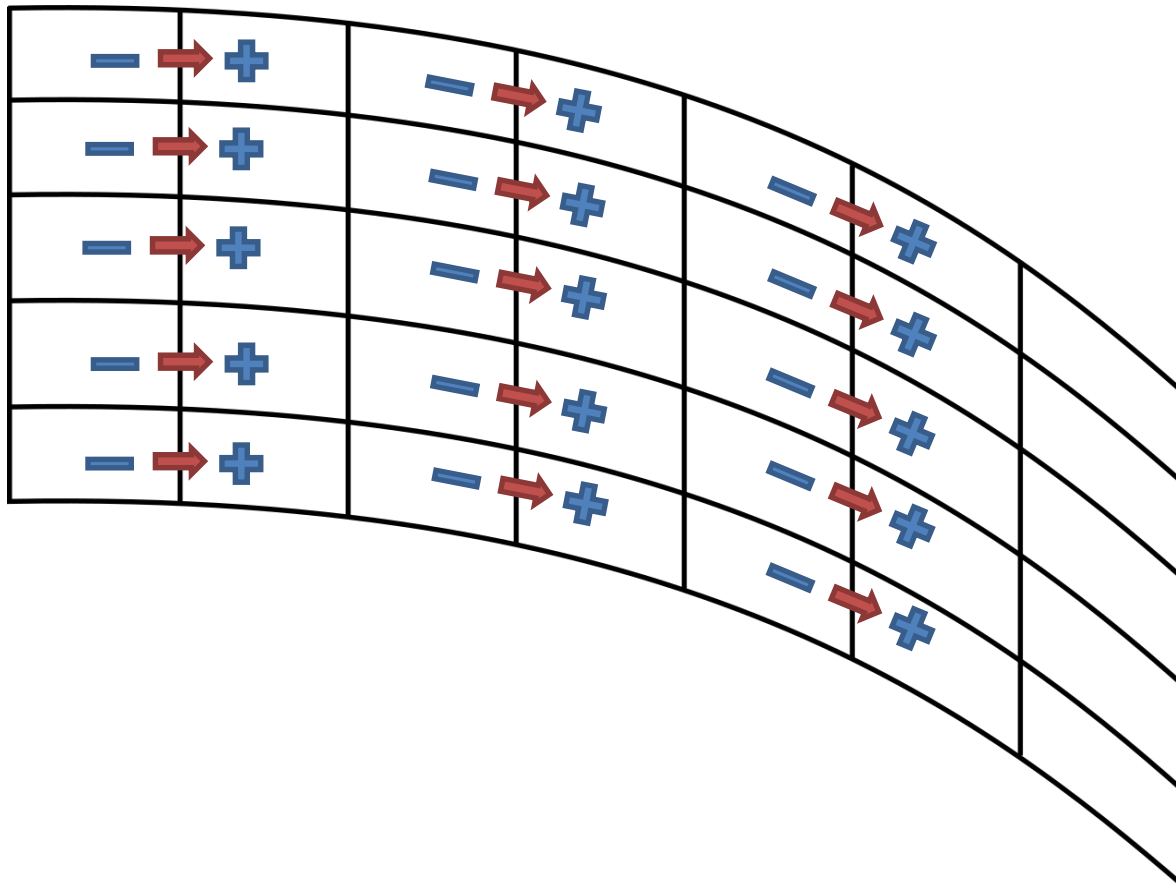
- Direction-based mapping thread safe and less redundancy



Convective Flux (2/3):

Thread mapping possibilities

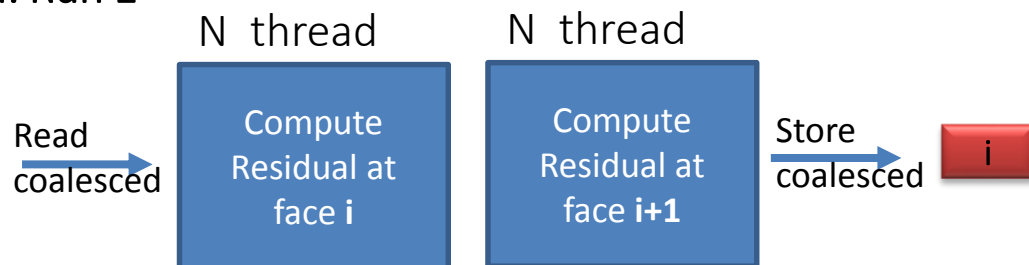
- Multicoloring (MC) Face-based mapping thread safe and No redundancy



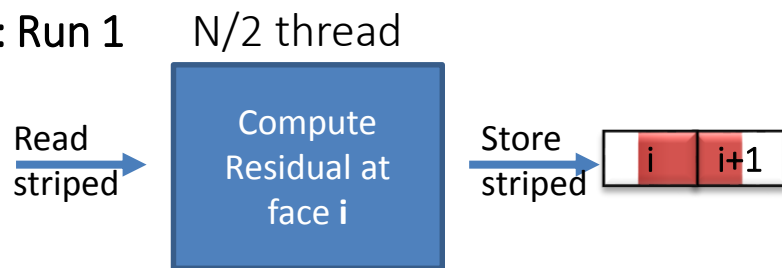
Convective Flux (3/3):

Multicolored (MC) vs redundant (Red)

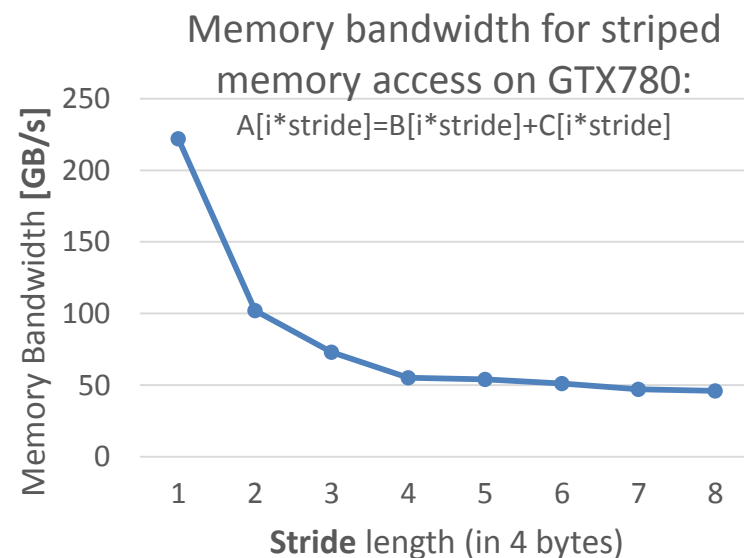
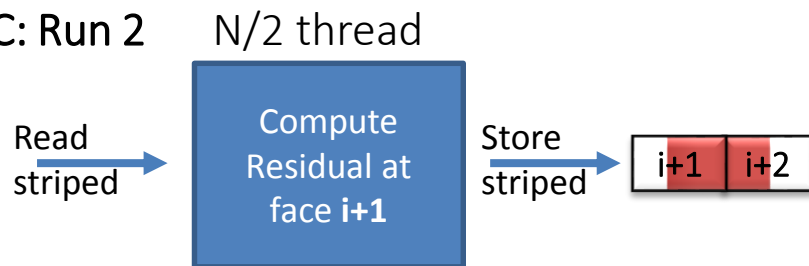
Red: Run 1



MC: Run 1

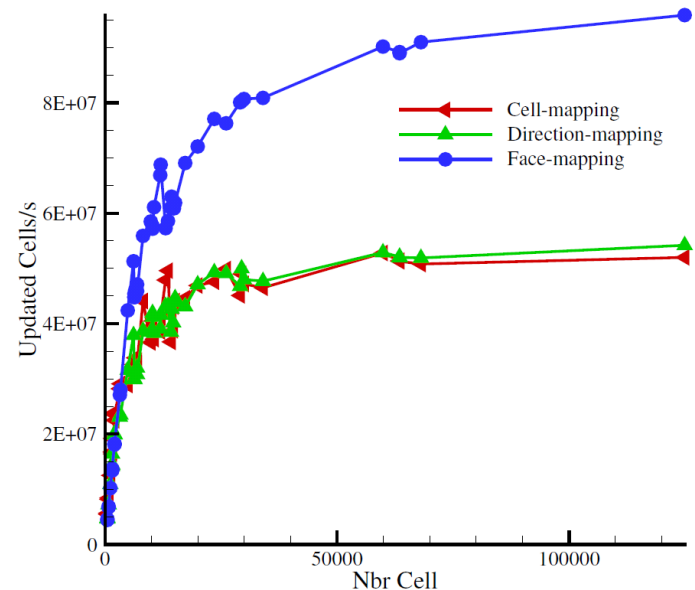


MC: Run 2



Convective Flux (3/3): Multicolored (MC) vs redundant (Red)

| | MC | RED |
|----------------------|-------|------|
| Face fluxes per call | $N/2$ | $2N$ |
| Total faces fluxes | N | $2N$ |
| Time per call [ms] | 0,28 | 0,71 |
| Total time [ms] | 0,56 | 0,71 |
| Operations ratio | - | 2x |
| Total Speedup | 1,26x | - |

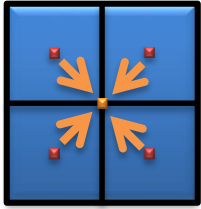


1,26x instead of 2x:
cost of striped access

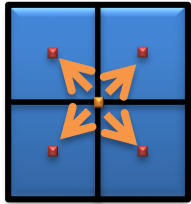
Convergence Acceleration on GPU (1/3)

- Explicit solver is well adapted to the GPU architecture
- Flow convergence is slow (CFL limitation)
- Need for convergence acceleration.
- convergence acceleration methods on the GPU?
 - Multigrid
 - Implicit residual smoothing

Convergence Acceleration on GPU (2/3): Multigrid is also fast on the GPU

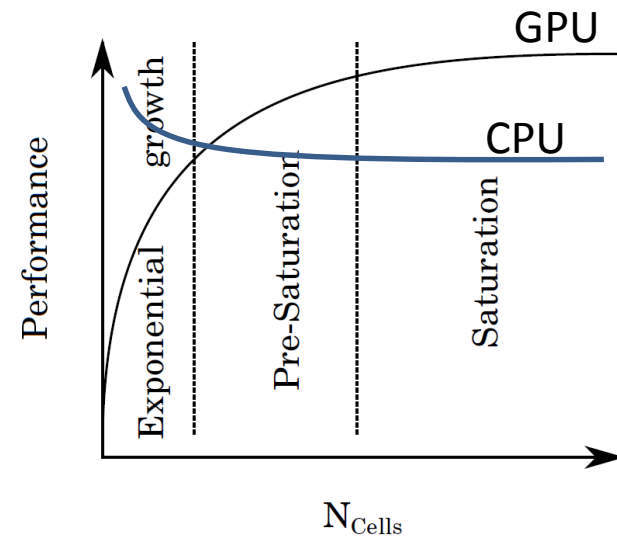
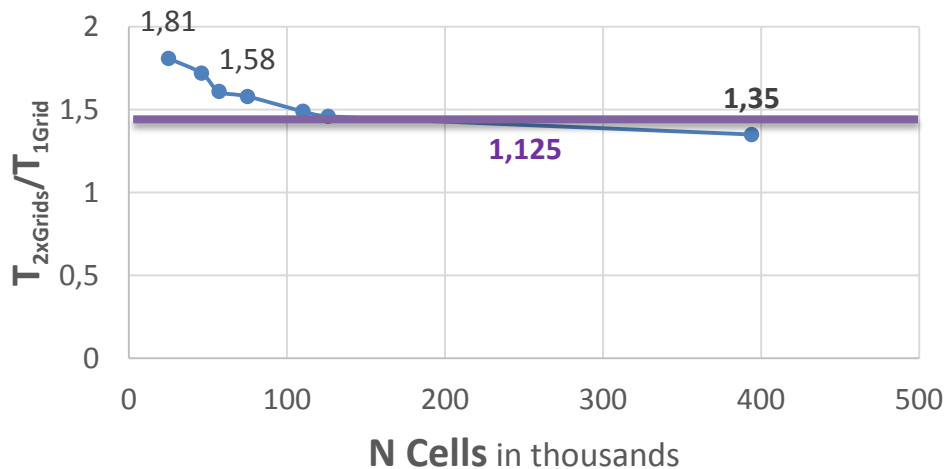


- Solve on fine grid
- Interpolate solution and residual to coarse grid



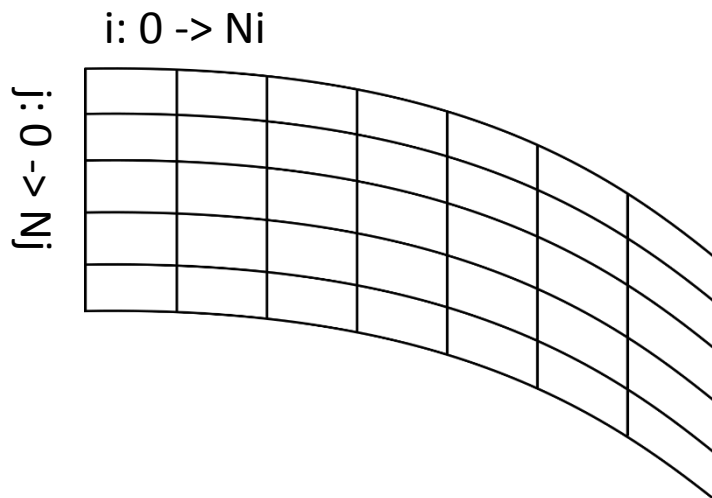
- Solve on coarse grid assisted by fine residual
- Prolongate coarse correction to fine grid

Cost of a 2-Grid scheme
converging to ideal cost of 1,125



Convergence Acceleration on GPU (3/3): Implicit Residual Smoothing on GPU

- Higher CFL \rightarrow Oscillation in the solution.
- A smoother residual reduces the oscillation \rightarrow Higher CFLs.
- Smoothing: diffusion equation \rightarrow solve a tridiagonal system.



$i: 0 \rightarrow 2 \ 3$

j

| | | | |
|----|----|----|----|
| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |
| 16 | 17 | 18 | 19 |

j

$j=0$

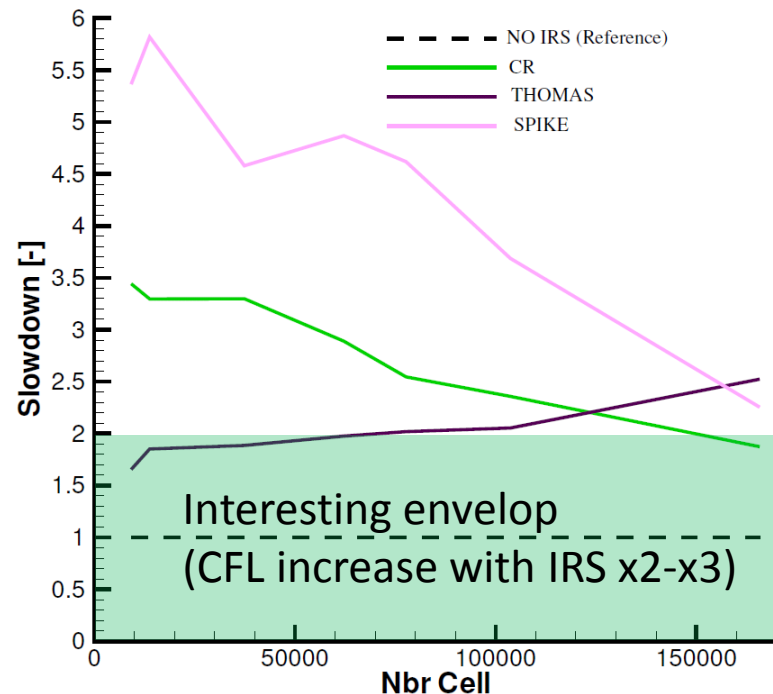
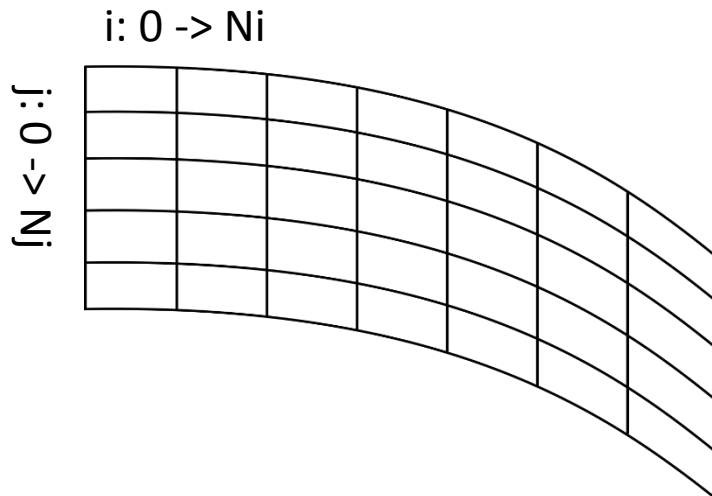
$$\begin{array}{c}
 \begin{array}{ccc|c}
 0 & (1+2\varepsilon_0) & -\varepsilon_0 & R_0^* \\
 -\varepsilon_1 & (1+2\varepsilon_1) & -\varepsilon_1 & R_1^* \\
 -\varepsilon_2 & (1+2\varepsilon_2) & -\varepsilon_2 & R_2^* \\
 0 & -\varepsilon_3 & (1+2\varepsilon_3) & R_3^*
 \end{array}
 \end{array}
 =
 \begin{array}{c}
 R_0 \\
 R_1 \\
 R_2 \\
 R_3
 \end{array}$$

$j=1$

$$\begin{array}{c}
 \begin{array}{ccc|c}
 0 & (1+2\varepsilon_4) & -\varepsilon_4 & R_4^* \\
 -\varepsilon_5 & (1+2\varepsilon_5) & -\varepsilon_5 & R_5^* \\
 -\varepsilon_6 & (1+2\varepsilon_6) & -\varepsilon_6 & R_6^* \\
 0 & -\varepsilon_7 & (1+2\varepsilon_7) & R_7^*
 \end{array}
 \end{array}
 =
 \begin{array}{c}
 R_4 \\
 R_5 \\
 R_6 \\
 R_7
 \end{array}$$

Convergence Acceleration on GPU (3/3): Implicit Residual Smoothing on GPU

- Higher CFL \rightarrow Oscillation in the solution.
- A smoother residual reduces the oscillation \rightarrow Higher CFLs.
- Smoothing: diffusion equation \rightarrow solve a tridiagonal system.



Main objective:

A more tangible GPU potential



CFD GPU solvers

- Explicit time integration
- Implicit time integration



Classification
of CFD operations



Proof-of-concept:
Optimization cases



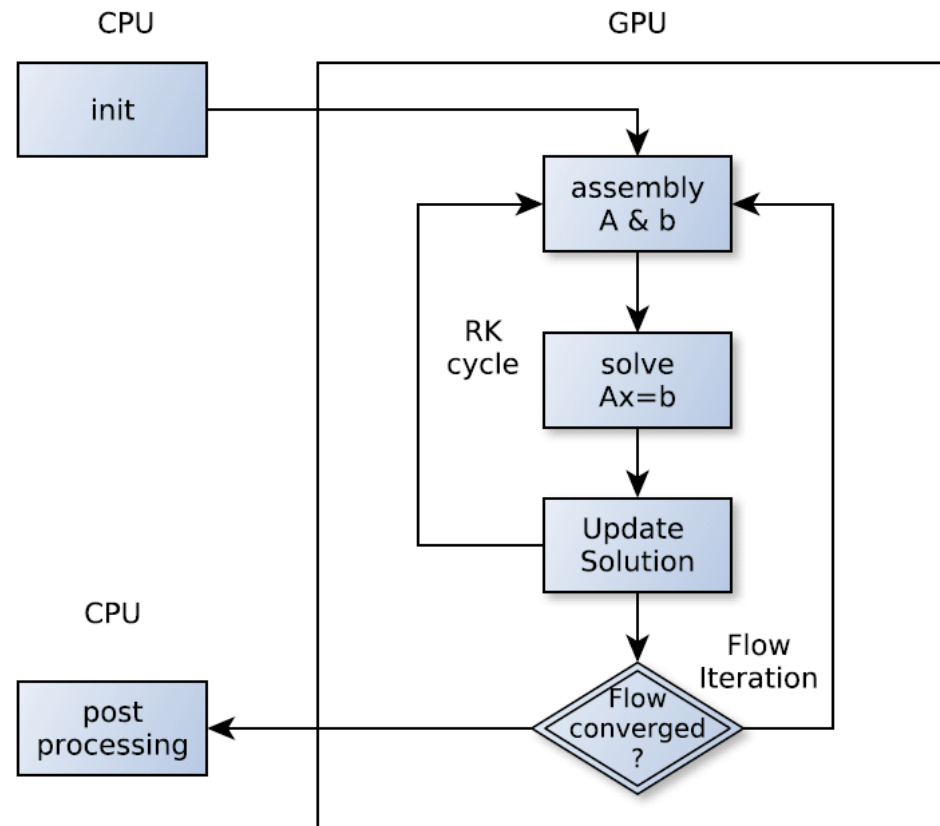
Summary
and Conclusions

Implicit Time Stepping is more Stable but ...

$$\frac{\Omega I}{\Delta t} \Delta \vec{W}^n = -\vec{R}^{(n+1)}$$

$$\vec{R}^{n+1} \approx \vec{R}^n + \left(\frac{\delta \vec{R}}{\delta \vec{W}} \right) \Delta \vec{W}^n$$

$$\left[\frac{\Omega I}{\Delta t} + \left(\frac{\delta \vec{R}}{\delta \vec{W}} \right) \right] \Delta \vec{W}^n = \vec{R}^n$$



GMRES + Preconditioner

$$Ax = b,$$

$$AM^{-1}u = b, x \equiv M^{-1}u$$

$$\left[\frac{\Omega I}{\Delta t} + \left(\frac{\delta \vec{R}}{\delta \vec{W}} \right) \right] \Delta \vec{W}^n = \vec{R}^n$$

Algorithm 4 preconditioned GMRES

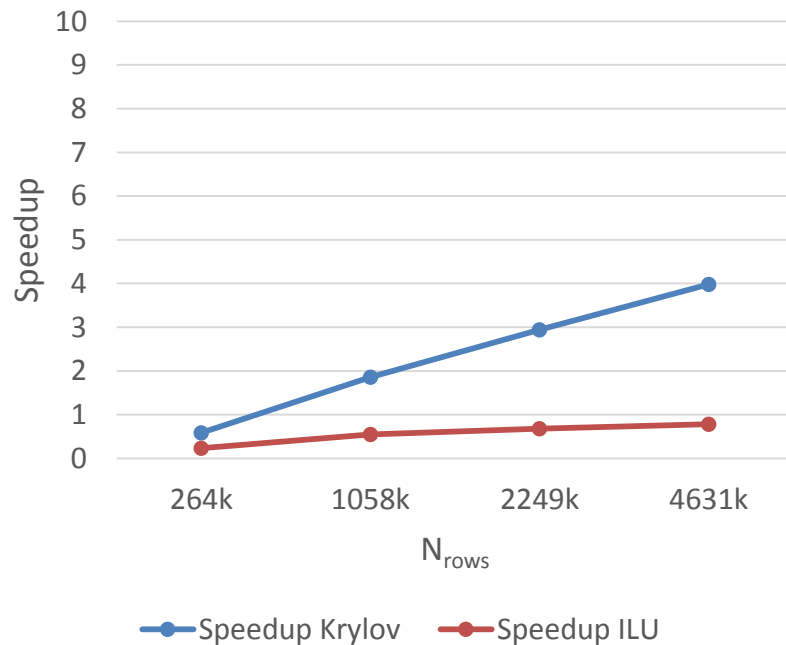
```
1:  $r_0 = b - Ax_0$ ,  $\beta := \|r_0\|_2$  and  $v_1 := r_0/\beta$ 
2: while  $\|r\|_2 > \epsilon \|b\|_2$  do
3:   for  $j=1$  to  $m$  do
4:      $w_j := AM^{-1}v_j$ 
5:     for  $i=1$  to  $j$  do
6:        $h_{ij} = (w_i, v_i)$ 
7:        $w_j := w_j - h_{ij}v_i$ 
8:     end for
9:      $h_{j+1,j} = \|w_j\|_2$  and  $v_{j+1} = w_j/h_{j+1,j}$ 
10:     $V_m := [v_1, \dots, v_m]$ ,  $\bar{H} = h_{ij, 1 \leq i < m+1, 1 \leq j < m}$ 
11:  end for
12:   $y_m = \operatorname{argmin}_y \|\beta e_1 - \bar{H}_m y\|$ 
13:   $x_m = x_0 + M^{-1}V_m y_m$ 
14:   $x_0 := x_m$ 
15: end while
```

Algorithm 5 ILU(0)

```
1: for  $i=2$  to  $n$  do
2:   for  $k=1$  to  $i-1$  do
3:     if  $(i, k) \in S$  then
4:        $a_{ik} = a_{ik}/a_{kk}$ 
5:       for  $j=k+1$  to  $n$  do
6:         if  $(i, j) \in S$  then
7:            $a_{ij} = a_{ik}/a_{kj}$ 
8:         end if
9:       end for
10:    end if
11:  end for
12: end for
```

ILU is costly on GPU

- **ILU-GMRES:** Small gain on every iteration but ILU setup is slow:

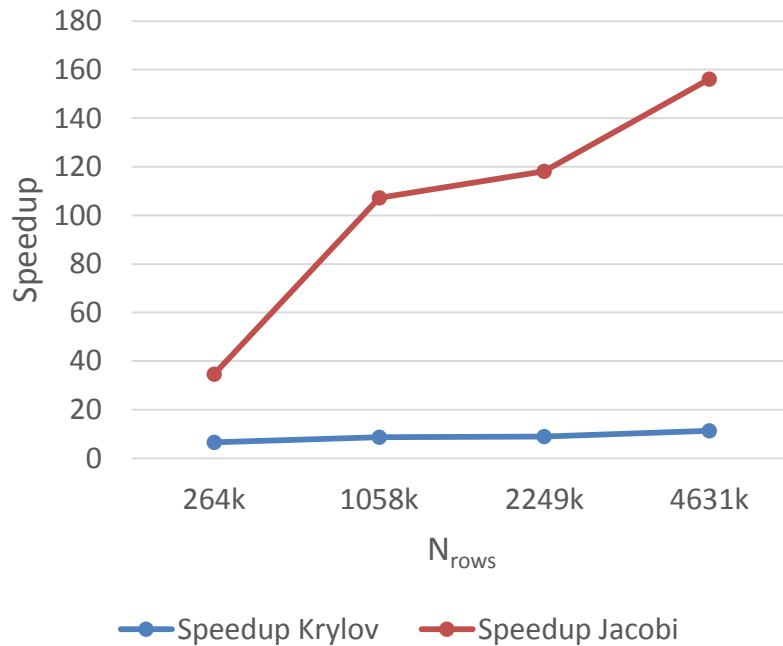


- **MCILU-GMRES:** Multi-colored ILU fast only for small problems.

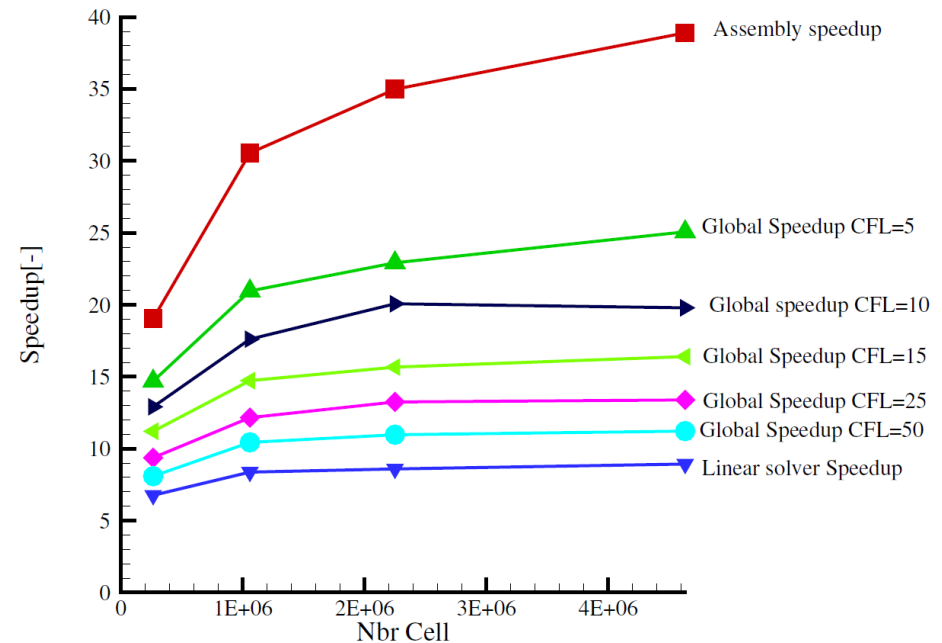


Why not Jacobi PC

- Jacobi-GMRES: very fast but stable only for small time steps

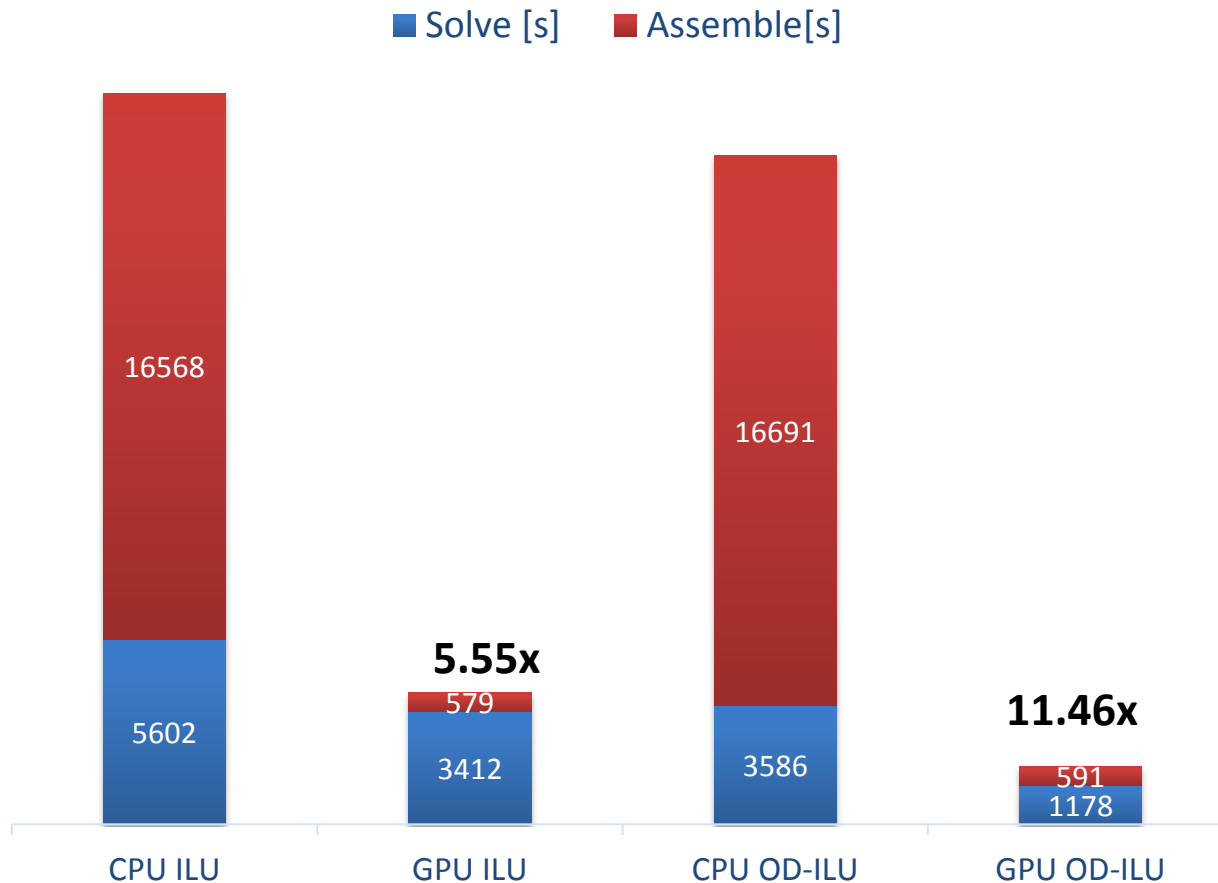


- Jacobi-GMRES: Speedup decreases for higher CFLs



On-demand factorization

```
if (itr > MAX_ITR ) M <- LU_Factorization (A)  
(x, itr) <- FGMRES (A,M,b)
```



Main objective:

A more tangible GPU potential



CFD GPU solvers



Classification
of CFD operations



Proof-of-concept:
Optimization cases

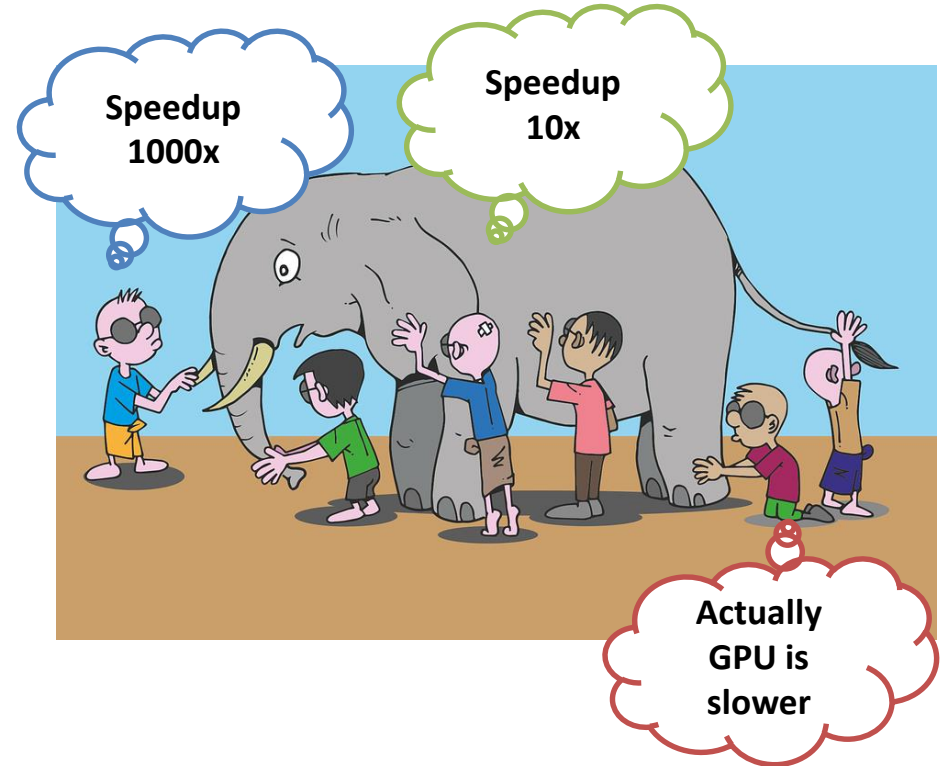


Summary
and Conclusions

Classification (1/2): GPUs controversy

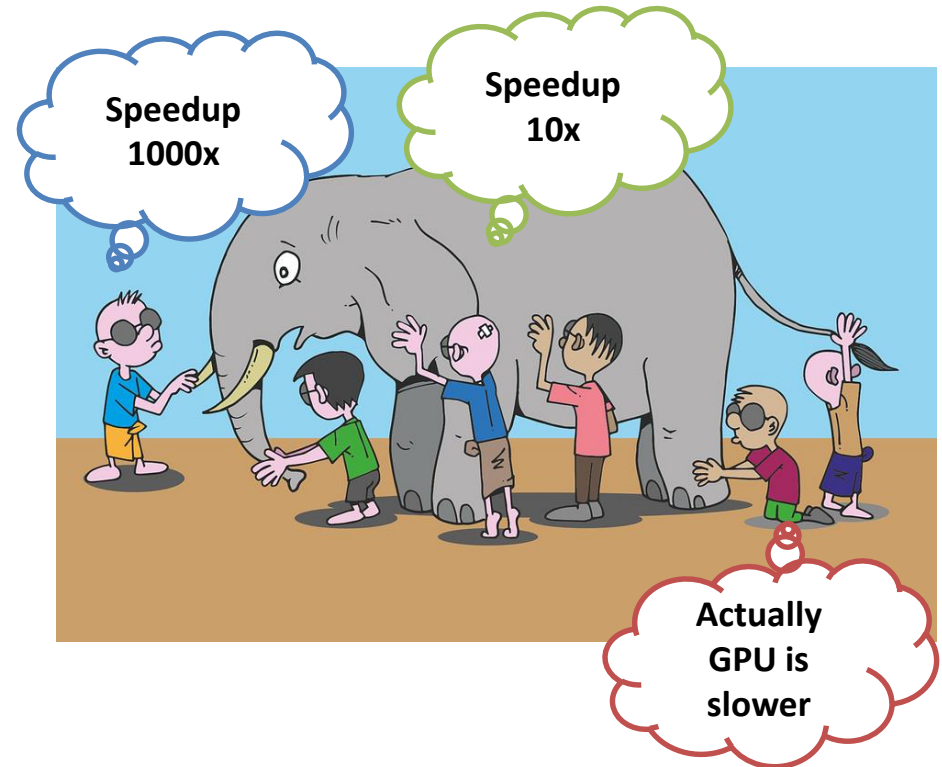
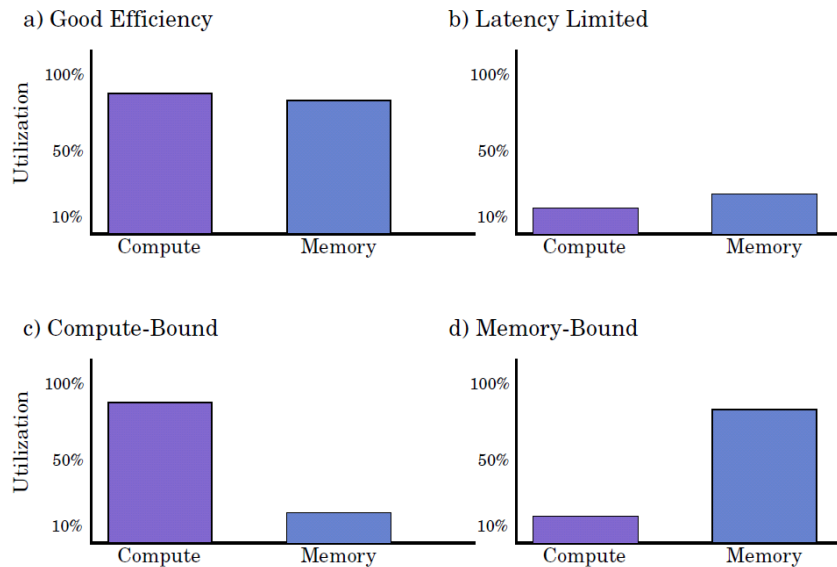
- GPU thousands of lightweight cores.
- Explicit solver: 10x to 100x speedup.
- Implicit solver: 1x to 10x speedup

→ We need a classification



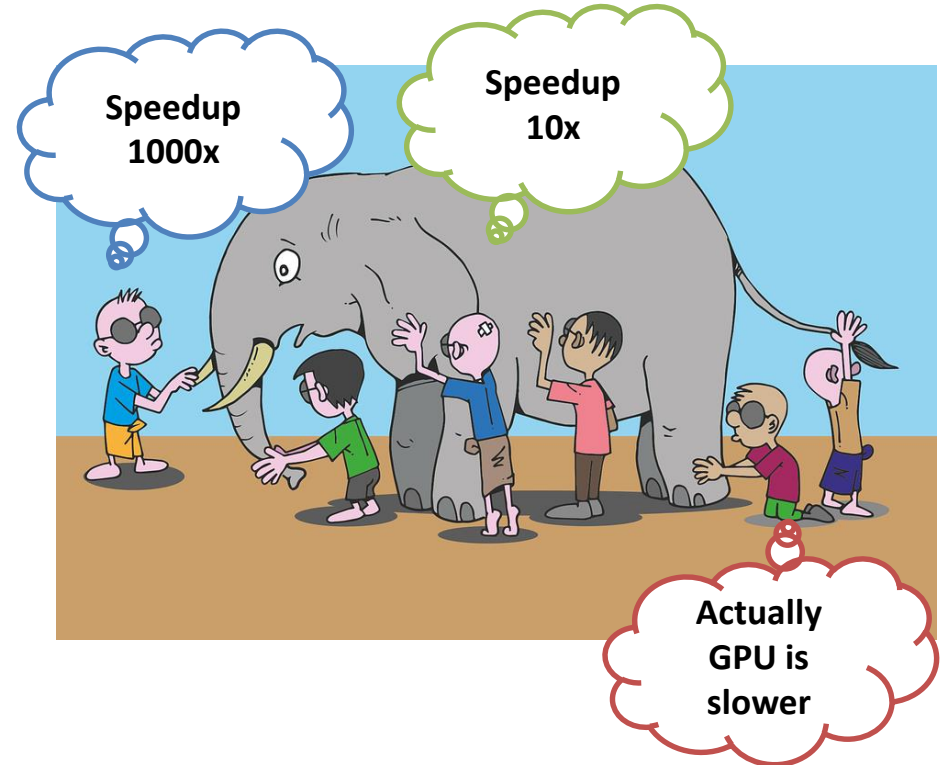
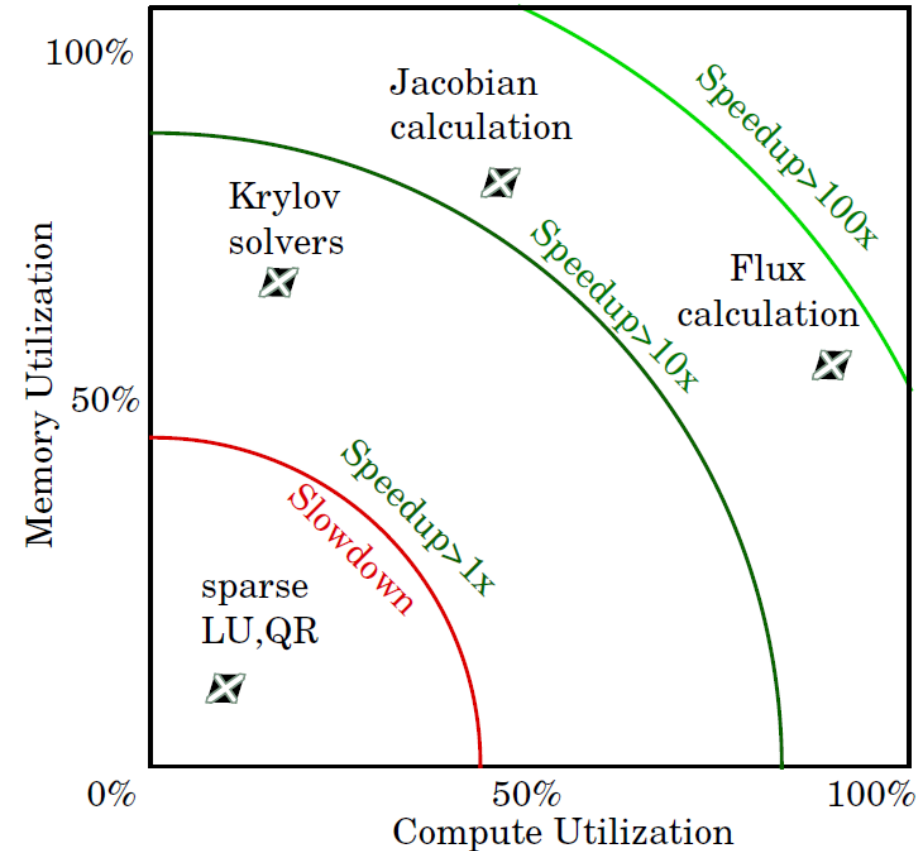
Picture modified from :<https://pixabay.com/en/blind-men-elephant-story-feel-see-1458438/>

Classification (1/2): GPUs controversy



Picture modified from :<https://pixabay.com/en/blind-men-elephant-story-feel-see-1458438/>

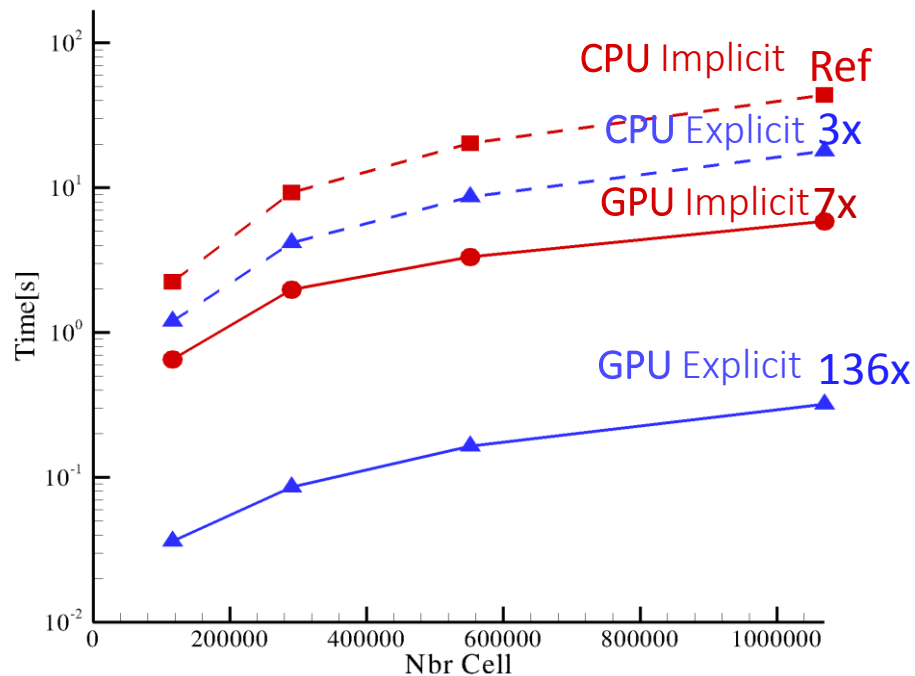
Classification (2/2): CFD operations



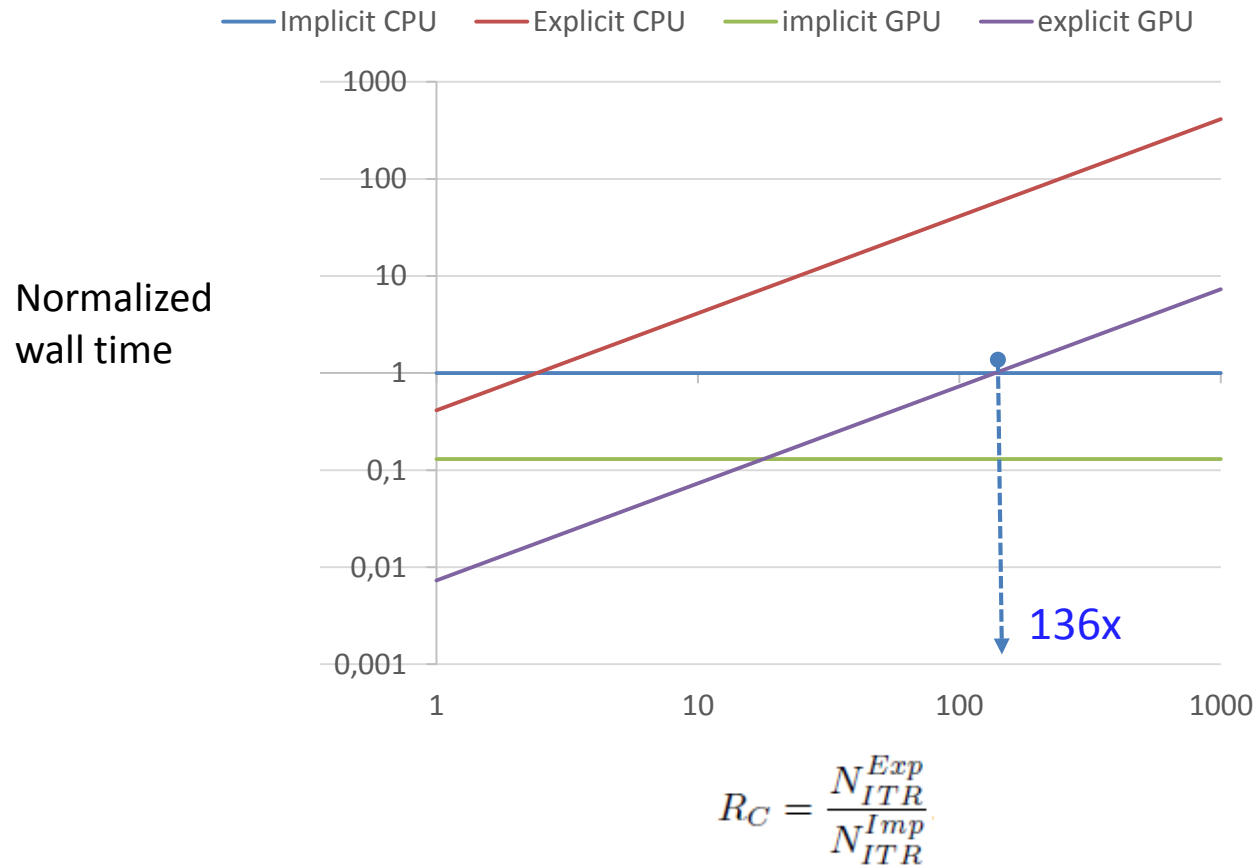
Full article:

[Aissa M., Verstraete, T., and Vuik, c.. "Toward a GPU-aware comparison of explicit and implicit CFD simulations on structured meshes." Computers & Mathematics with Applications 74.1 \(2017\): 201-217.](#)

Performance Comparison: Explicit/Implicit



Performance Comparison: Explicit/Implicit



Main objective:

A more tangible GPU potential



CFD GPU solvers



Classification
of CFD operations



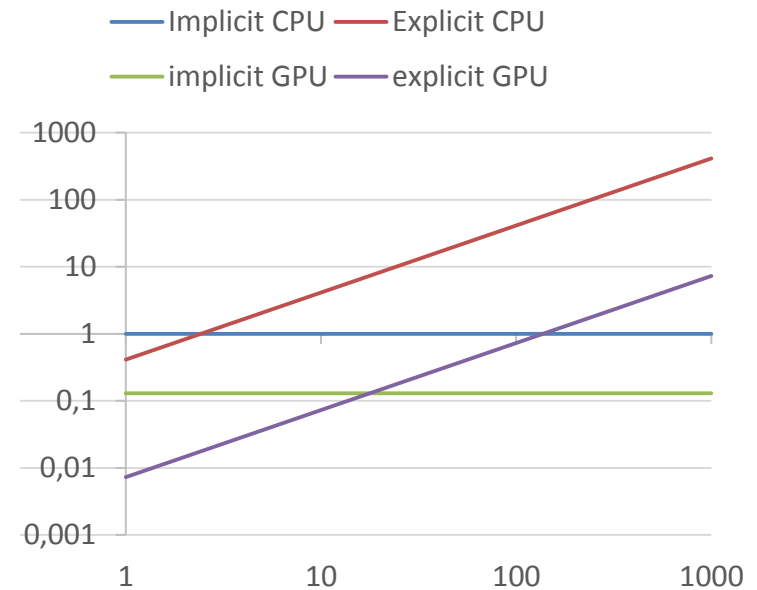
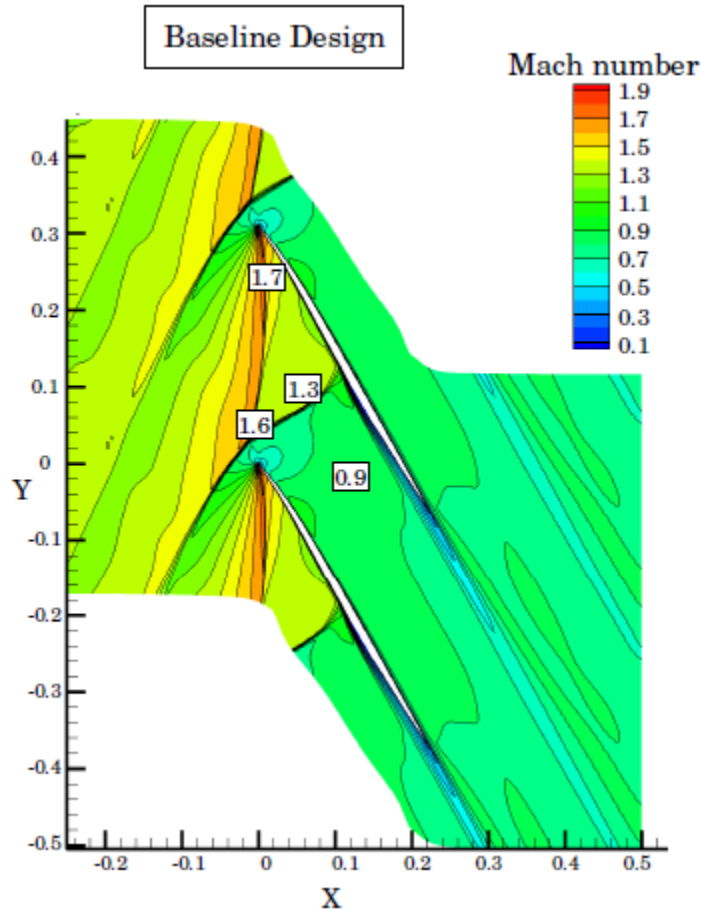
Proof-of-concept:
Optimization cases



Summary
and Conclusions

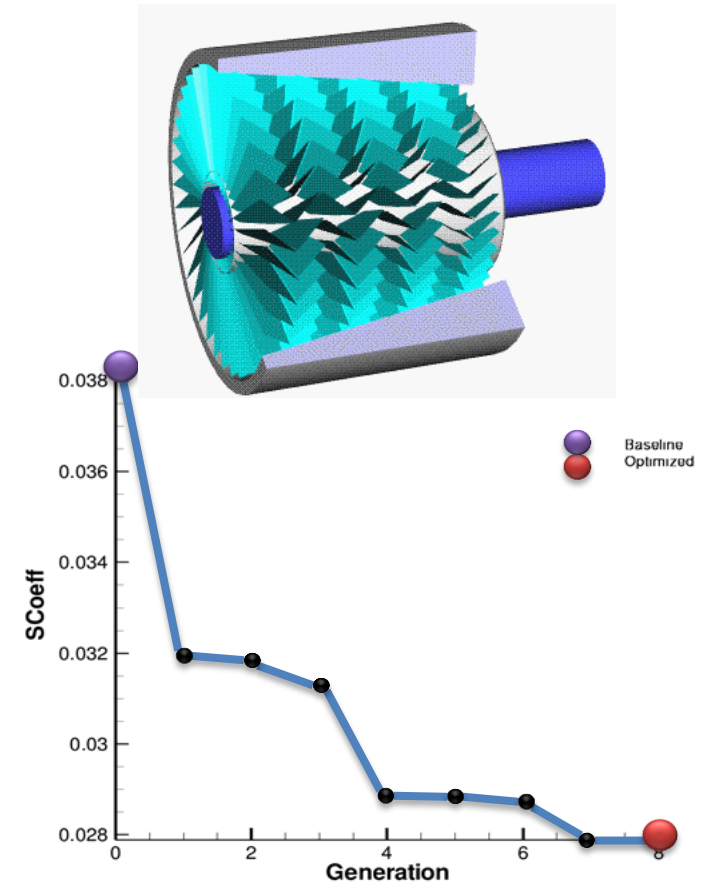
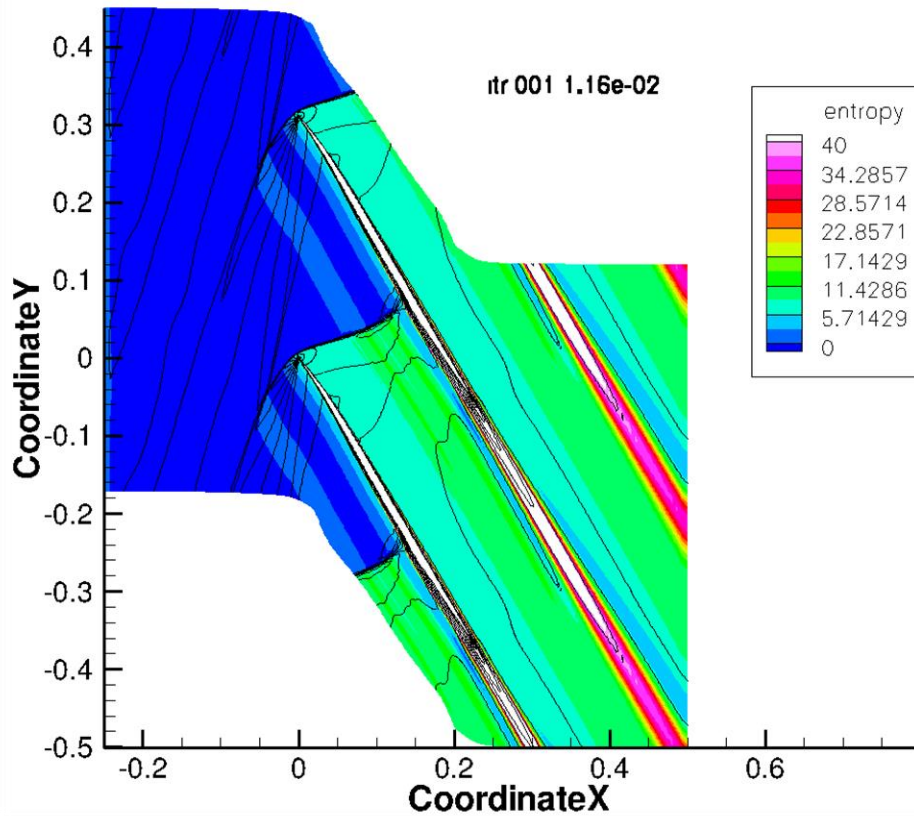
Example of a stator Optimization

$R_c=14$ (low CFL for implicit =15)

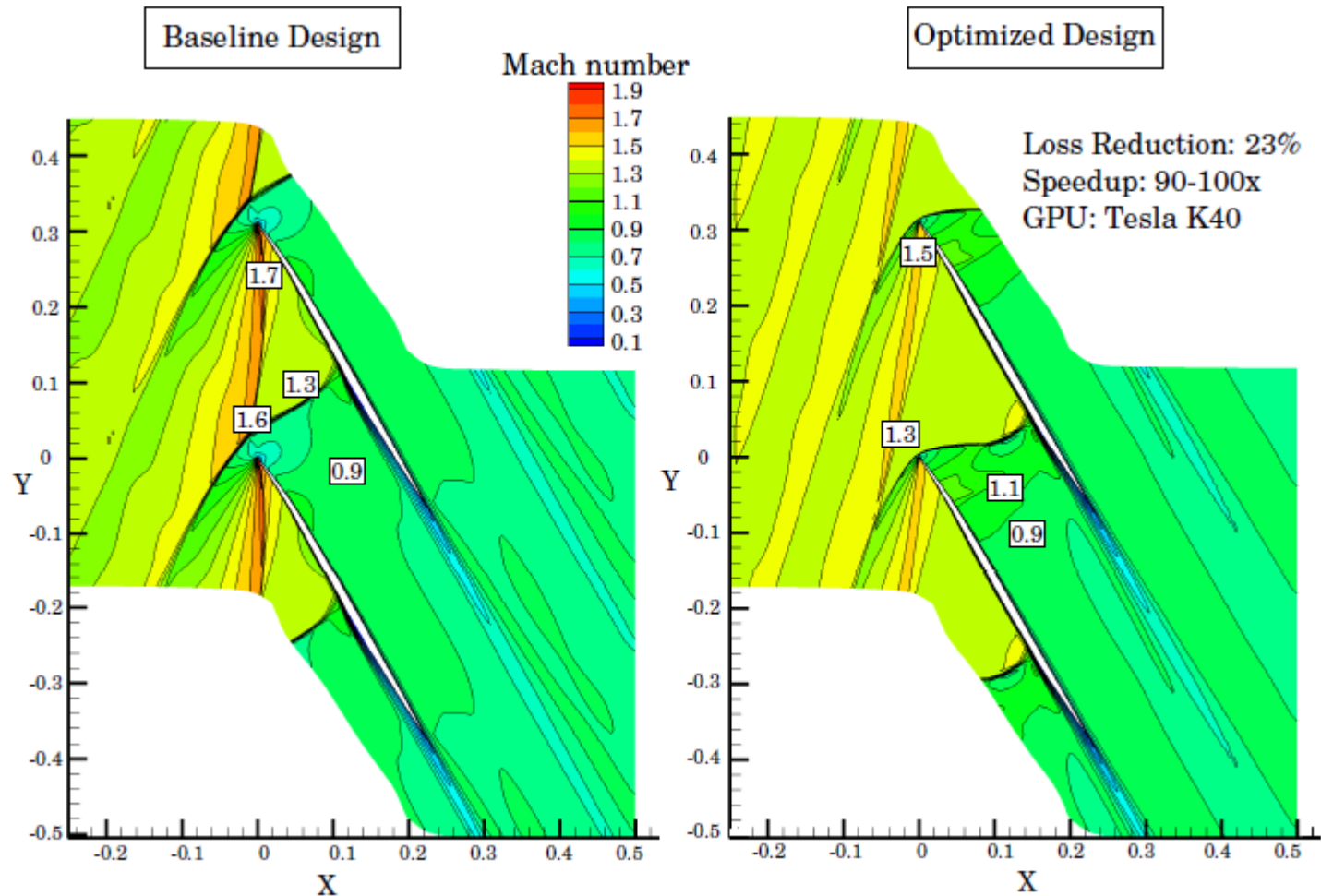


$$R_C = \frac{N_{ITR}^{Exp}}{N_{ITR}^{Imp}}$$

Example of a stator Optimization

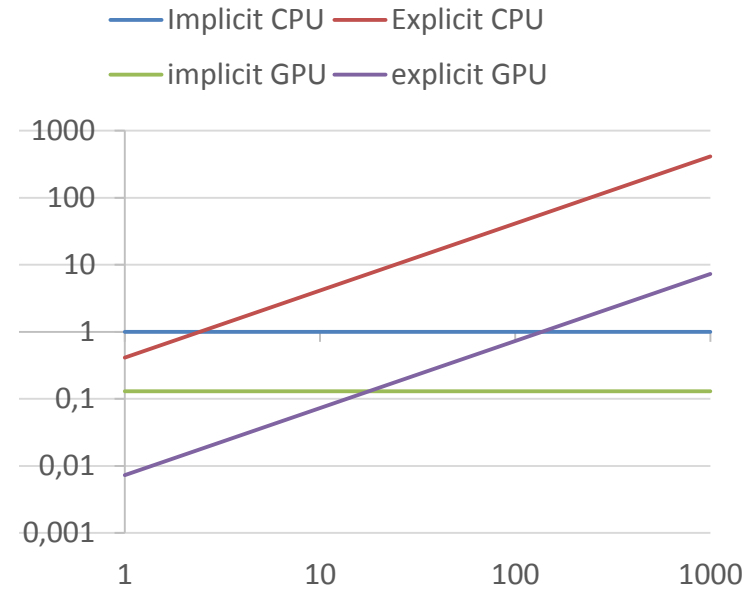
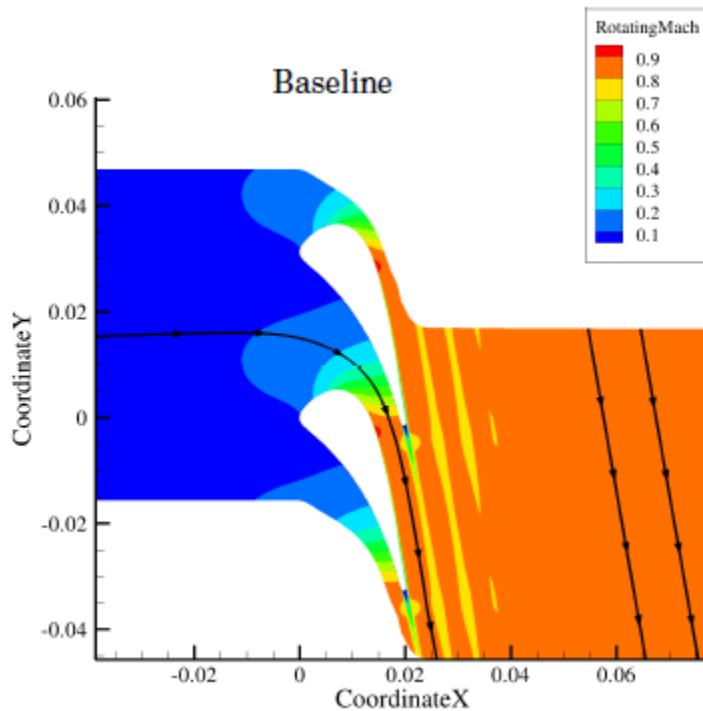


Example of a stator Optimization



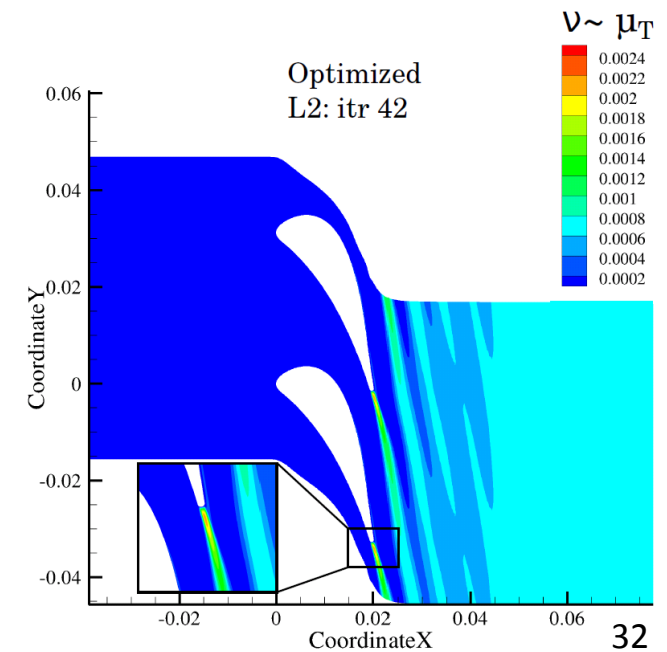
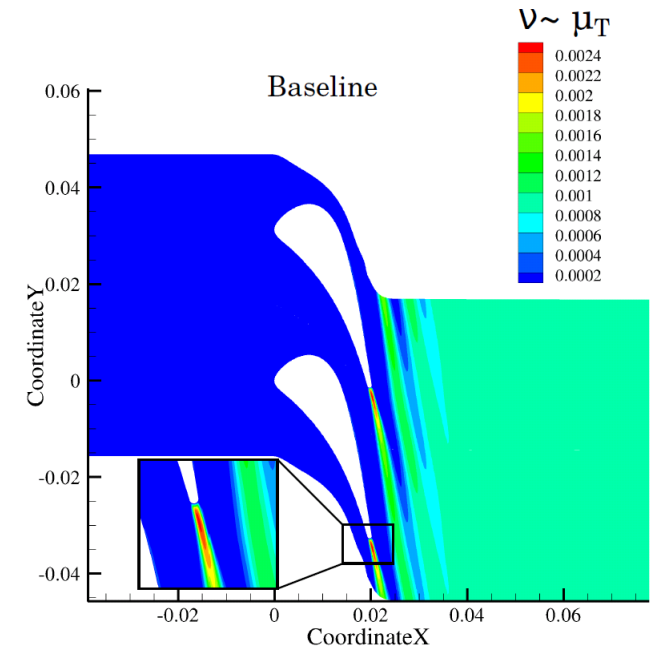
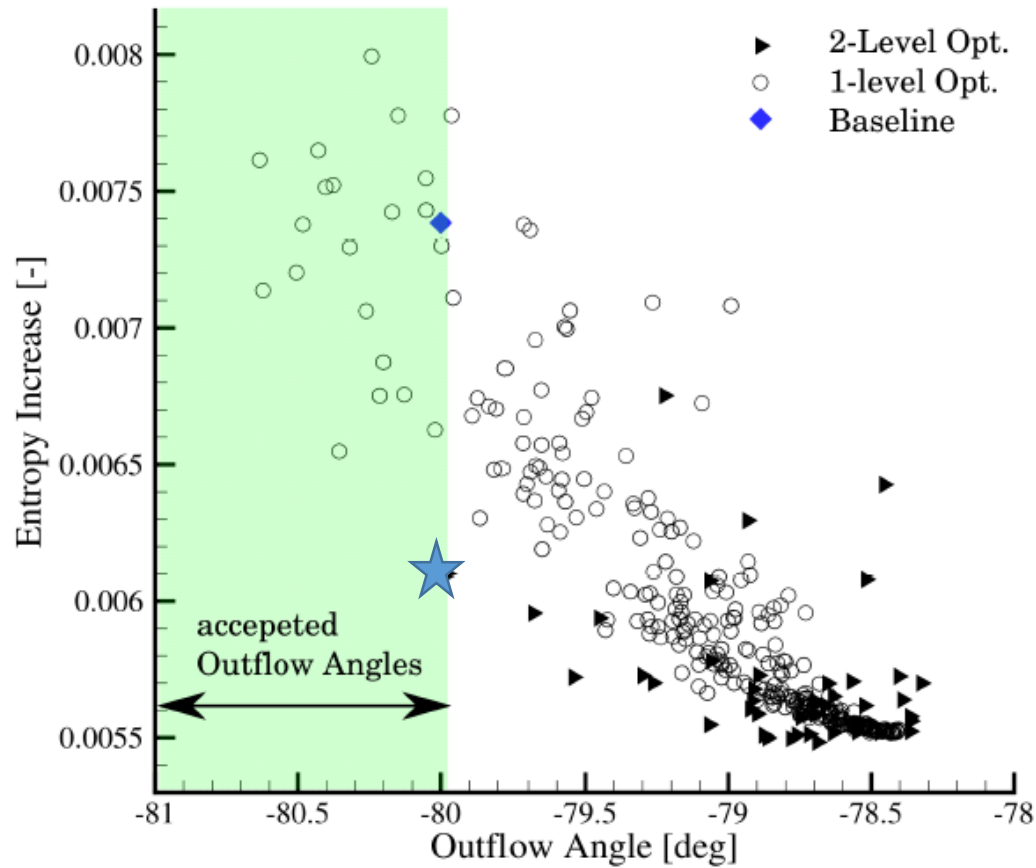
LS82 cascade

$R_c=457$ (Explicit solver bad flow convergence)

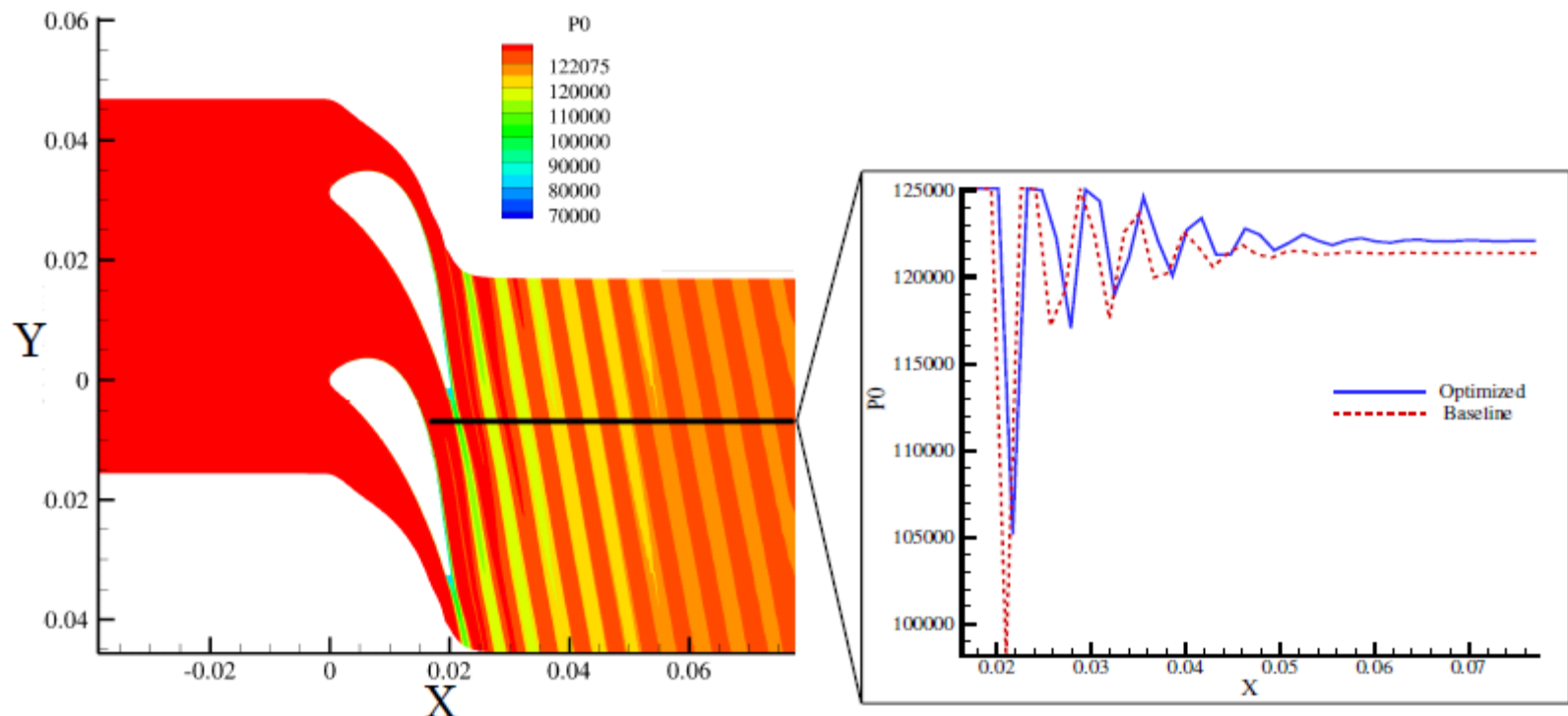


$$R_C = \frac{N_{ITR}^{Exp}}{N_{ITR}^{Imp}}$$

LS82 cascade: Results



LS82 cascade: Optimized blade



Main objective:

A more tangible GPU potential



CFD GPU solvers



Classification
of CFD operations



Proof-of-concept:
Optimization cases



Summary
and Conclusions

Summary



Explicit RANS:
100x-180x speedup.



Implicit RANS: 10x-20x speedup
(due to slow preconditioning).



On-demand preconditioning:
x3 faster but GPU-friendlier
preconditioner is needed.



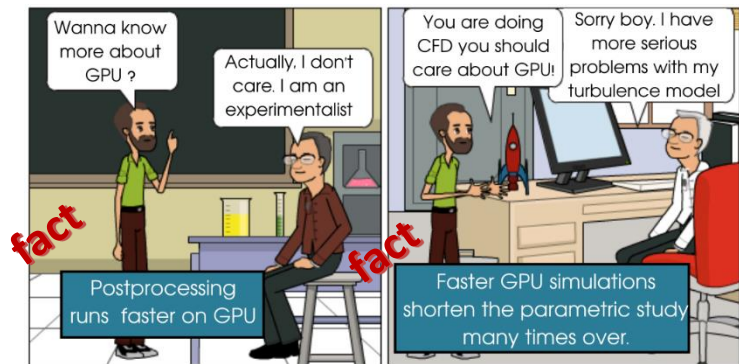
The classification:
an operation-specific
acceleration offers more insights.



Choice Explicit/Implicit:
Convergence ratio is decisive.

Can your simulation profit from the GPU?

- Where you situate your algorithm (slide 4: QR to ray-tracing)?
- Do you need double precision (for half-precision FPGA is faster)?
- ready to code (otherwise openACC is easier to use)?
- Anyone provided a classification for operation used in your field?



Thanks for your attention

Dr. Mohamed H. Aissa

Turbomachinery & Propulsion Department

Email: aissa@vki.ac.be

www.researchgate.net/profile/Mohamed_Aissa3